# Amazon Fine Food Reviews Analysis

Data Source: https://www.kaggle.com/snap/amazon-fine-food-reviews

EDA: https://nycdatascience.com/blog/student-works/amazon-fine-foods-visualization/

The Amazon Fine Food Reviews dataset consists of reviews of fine foods from Amazon.

Number of reviews: 568,454
Number of users: 256,059
Number of products: 74,258
Timespan: Oct 1999 - Oct 2012
Number of Attributes/Columns in data: 10

Attribute Information:

1. Id
2. ProductId - unique identifier for the product
3. UserId - unqiue identifier for the user
4. ProfileName
5. HelpfulnessNumerator - number of users who found the review helpful
6. HelpfulnessDenominator - number of users who indicated whether they found the review helpful or not
7. Score - rating between 1 and 5
8. Time - timestamp for the review
9. Summary - brief summary of the review
10. Text - text of the review

**Objective:**

Given a review, determine whether the review is positive (rating of 4 or 5) or negative (rating of 1 or 2).

[Q] How to determine if a review is positive or negative?

[Ans] We could use Score/Rating. A rating of 4 or 5 can be cosnidered as a positive review. A rating of 1 or 2 can be considered as negative one. A review of rating 3 is considered nuetral and such reviews are ignored from our analysis. This is an approximate and proxy way of determining the polarity (positivity/negativity) of a review.

# [1]. Reading Data

```
In [0]:  from google.colab import drive
         drive.mount('/content/drive')
```

```
Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?
client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleuser
content.com&redirect_uri=urn%3Aietf%3Awg%3Aoauth%3A2.0%3Aoob&scope=emai
l%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdocs.test%20https%3A%2F%2
Fwww.googleapis.com%2Fauth%2Fdrive%20https%3A%2F%2Fwww.googleapis.com%2
Fauth%2Fdrive.photos.readonly%20https%3A%2F%2Fwww.googleapis.com%2Faut
h%2Fpeopleapi.readonly&response_type=code

Enter your authorization code:
..........
Mounted at /content/drive
```

## [1.1] Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

In order to load the data, We have used the SQLITE dataset as it is easier to query the data and visualise the data efficiently.

Here as we only want to get the global sentiment of the recommendations (positive or negative), we will purposefully ignore all Scores equal to 3. If the score is above 3, then the recommendation wil be set to "positive". Otherwise, it will be set to "negative".

In [0]:
```python
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")


import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer

import re
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
import string
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer

from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle

from tqdm import tqdm
import os
```

```
/usr/local/lib/python3.6/dist-packages/smart_open/ssh.py:34: UserWarnin
g: paramiko missing, opening SSH/SCP/SFTP paths will be disabled.  `pip
install paramiko` to suppress
  warnings.warn('paramiko missing, opening SSH/SCP/SFTP paths will be d
isabled.  `pip install paramiko` to suppress')
```

In [0]:
```
con = sqlite3.connect('/content/drive/My Drive/Colab Notebooks/databas
e.sqlite')
filtered_data_100k = pd.read_sql_query("SELECT * FROM Reviews WHERE Sco
re != 3 LIMIT 100000", con)
def partition(x):
    if x < 3:
        return 0
    return 1

actualScore = filtered_data_100k['Score']
positiveNegative = actualScore.map(partition)
filtered_data_100k['Score'] = positiveNegative
print("Number of data points in our data", filtered_data_100k.shape)
filtered_data_100k.head(1)
```

Number of data points in our data (100000, 10)

Out[0]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfulnes |
|---|---|---|---|---|---|---|
| 0 | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1 | 1 |

In [0]:
```
con = sqlite3.connect('/content/drive/My Drive/Colab Notebooks/databas
e.sqlite')
filtered_data_20k = pd.read_sql_query("SELECT * FROM Reviews WHERE Scor
```

```
e != 3 LIMIT 20000", con)
def partition(x):
    if x < 3:
        return 0
    return 1

actualScore = filtered_data_20k['Score']
positiveNegative = actualScore.map(partition)
filtered_data_20k['Score'] = positiveNegative
print("Number of data points in our data", filtered_data_20k.shape)
filtered_data_20k.head(1)
```

Number of data points in our data (20000, 10)

Out[0]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfulnes |
|---|---|---|---|---|---|---|
| 0 | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1 | 1 |

In [0]:
```
display = pd.read_sql_query("""
SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
FROM Reviews
GROUP BY UserId
HAVING COUNT(*)>1
""", con)
```

In [0]:
```
print(display.shape)
display.head()
```

(80668, 7)

Out[0]:

| | UserId | ProductId | ProfileName | Time | Score | Text | COUI |
|---|---|---|---|---|---|---|---|
| 0 | #oc-R115TNMSPFT9I7 | B007Y59HVM | Breyton | 1331510400 | 2 | Overall its just OK when considering the price... | 2 |
| 1 | #oc-R11D9D7SHXIJB9 | B005HG9ET0 | Louis E. Emory "hoppy" | 1342396800 | 5 | My wife has recurring extreme muscle spasms, u... | 3 |
| 2 | #oc-R11DNU2NBKQ23Z | B007Y59HVM | Kim Cieszykowski | 1348531200 | 1 | This coffee is horrible and unfortunately not ... | 2 |
| 3 | #oc-R11O5J5ZVQE25C | B005HG9ET0 | Penguin Chick | 1346889600 | 5 | This will be the bottle that you grab from the... | 3 |
| 4 | #oc-R12KPBODL2B5ZD | B007OSBE1U | Christopher P. Presta | 1348617600 | 1 | I didnt like this coffee. Instead of telling y... | 2 |

```
In [0]: display[display['UserId']=='AZY10LLTJ71NX']
```

Out[0]:

| | UserId | ProductId | ProfileName | Time | Score | Text | C |
|---|---|---|---|---|---|---|---|

| | UserId | ProductId | ProfileName | Time | Score | Text | |
|---|---|---|---|---|---|---|---|
| **80638** | AZY10LLTJ71NX | B006P7E5ZI | undertheshrine "undertheshrine" | 1334707200 | 5 | I was recommended to try green tea extract to ... | 5 |

In [0]: `display['COUNT(*)'].sum()`

Out[0]: 393063

# [2] Exploratory Data Analysis

## [2.1] Data Cleaning: Deduplication

It is observed (as shown in the table below) that the reviews data had many duplicate entries. Hence it was necessary to remove duplicates in order to get unbiased results for the analysis of the data. Following is an example:

In [0]:
```
display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID
""", con)
display.head()
```

Out[0]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfuln |
|---|---|---|---|---|---|---|

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfulr |
|---|---|---|---|---|---|---|
| 0 | 78445 | B000HDL1RQ | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |
| 1 | 138317 | B000HDOPYC | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |
| 2 | 138277 | B000HDOPYM | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |
| 3 | 73791 | B000HDOPZG | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |
| 4 | 155049 | B000PAQ75C | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |

As it can be seen above that same user has multiple reviews with same values for HelpfulnessNumerator, HelpfulnessDenominator, Score, Time, Summary and Text and on doing analysis it was found that

ProductId=B000HDOPZG was Loacker Quadratini Vanilla Wafer Cookies, 8.82-Ounce Packages (Pack of 8)

ProductId=B000HDL1RQ was Loacker Quadratini Lemon Wafer Cookies, 8.82-Ounce Packages (Pack of 8) and so on

It was inferred after analysis that reviews with same parameters other than ProductId belonged to the same product just having different flavour or quantity. Hence in order to reduce redundancy it was decided to eliminate the rows having same parameters.

The method used for the same was that we first sort the data according to ProductId and then just keep the first similar product review and delelte the others. for eg. in the above just the review for ProductId=B000HDL1RQ remains. This method ensures that there is only one representative for each product and deduplication without sorting would lead to possibility of different representatives still existing for the same product.

```python
In [0]: #Sorting data according to ProductId in ascending order
        sorted_data_100k=filtered_data_100k.sort_values('ProductId', axis=0, as
        cending=True, inplace=False, kind='quicksort', na_position='last')
```

```python
In [0]: #Sorting data according to ProductId in ascending order
        sorted_data_20k=filtered_data_20k.sort_values('ProductId', axis=0, asce
        nding=True, inplace=False, kind='quicksort', na_position='last')
```

```python
In [0]: #Deduplication of entries
        final_100k=sorted_data_100k.drop_duplicates(subset={"UserId","ProfileNa
        me","Time","Text"}, keep='first', inplace=False)
        final_100k.shape
```

```
Out[0]: (87775, 10)
```

```python
In [0]: #Deduplication of entries
        final_20k=sorted_data_20k.drop_duplicates(subset={"UserId","ProfileNam
        e","Time","Text"}, keep='first', inplace=False)
        final_20k.shape
```

```
Out[0]: (19354, 10)
```

```
In [0]:  #Checking to see how much % of data still remains
         (final_100k['Id'].size*1.0)/(filtered_data_100k['Id'].size*1.0)*100
```

```
Out[0]: 87.775
```

```
In [0]:  #Checking to see how much % of data still remains
         (final_20k['Id'].size*1.0)/(filtered_data_20k['Id'].size*1.0)*100
```

```
Out[0]: 96.77
```

**Observation:-** It was also seen that in two rows given below the value of HelpfulnessNumerator is greater than HelpfulnessDenominator which is not practically possible hence these two rows too are removed from calcualtions

```
In [0]:  display= pd.read_sql_query("""
         SELECT *
         FROM Reviews
         WHERE Score != 3 AND Id=44737 OR Id=64422
         ORDER BY ProductID
         """, con)

         display.head()
```

Out[0]:

|   | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfuln |
|---|------|-----------|----------------|-----------------------|----------------------|----------|
| 0 | 64422 | B000MIDROQ | A161DK06JJMCYF | J. E. Stephens "Jeanne" | 3 | 1 |

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfuln |
|---|---|---|---|---|---|---|
| 1 | 44737 | B001EQ55RW | A2V0I904FH7ABY | Ram | 3 | 2 |

```
In [0]: final_100k=final_100k[final_100k.HelpfulnessNumerator<=final_100k.Helpf
        ulnessDenominator]
```

```
In [0]: final_20k=final_20k[final_20k.HelpfulnessNumerator<=final_20k.Helpfulne
        ssDenominator]
```

```
In [0]: #Before starting the next phase of preprocessing lets see the number of
         entries left
        print(final_100k.shape)

        #How many positive and negative reviews are present in our dataset?
        final_100k['Score'].value_counts()
```

```
(87773, 10)
```

```
Out[0]: 1    73592
        0    14181
        Name: Score, dtype: int64
```

```
In [0]: #Before starting the next phase of preprocessing lets see the number of
         entries left
        print(final_20k.shape)

        #How many positive and negative reviews are present in our dataset?
        final_20k['Score'].value_counts()
```

```
(19354, 10)
```

```
Out[0]: 1    16339
        0     3015
        Name: Score, dtype: int64
```

# [3] Preprocessing

## [3.1]. Preprocessing Review Text

Now that we have finished deduplication our data requires some preprocessing before we go on further with analysis and making the prediction model.

Hence in the Preprocessing phase we do the following in the order below:-

1. Begin by removing the html tags
2. Remove any punctuations or limited set of special characters like , or . or # etc.
3. Check if the word is made up of english letters and is not alpha-numeric
4. Check to see if the length of the word is greater than 2 (as it was researched that there is no adjective in 2-letters)
5. Convert the word to lowercase
6. Remove Stopwords
7. Finally Snowball Stemming the word (it was obsereved to be better than Porter Stemming)

After which we collect the words used to describe positive and negative reviews

```python
In [0]:  # https://stackoverflow.com/a/47091490/4084039
         import re

         def decontracted(phrase):
             # specific
             phrase = re.sub(r"won't", "will not", phrase)
             phrase = re.sub(r"can\'t", "can not", phrase)

             # general
```

```python
        phrase = re.sub(r"n\'t", " not", phrase)
        phrase = re.sub(r"\'re", " are", phrase)
        phrase = re.sub(r"\'s", " is", phrase)
        phrase = re.sub(r"\'d", " would", phrase)
        phrase = re.sub(r"\'ll", " will", phrase)
        phrase = re.sub(r"\'t", " not", phrase)
        phrase = re.sub(r"\'ve", " have", phrase)
        phrase = re.sub(r"\'m", " am", phrase)
        return phrase
```

In [0]:
```python
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
# <br /><br /> ==> after the above steps, we are getting "br br"
# we are including them into stop words list
# instead of <br /> if we have <br/> these tags would have revmoved in
 the 1st step

stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've",\
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their',\
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after',\
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further',\
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more',\
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should',
```

```python
              "should've", 'now', 'd', 'll', 'm', 'o', 're', \
                've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't",
   'didn', "didn't", 'doesn', "doesn't", 'hadn',\
               "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "is
   n't", 'ma', 'mightn', "mightn't", 'mustn',\
               "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn',
    "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
               'won', "won't", 'wouldn', "wouldn't"])
```

In [0]:
```python
# Combining all the above stundents
from tqdm import tqdm
from bs4 import BeautifulSoup
preprocessed_reviews_100k = []
# tqdm is for printing the status bar
for sentance in tqdm(final_100k['Text'].values):
    sentance = re.sub(r"http\S+", "", sentance)
    sentance = BeautifulSoup(sentance, 'lxml').get_text()
    sentance = decontracted(sentance)
    sentance = re.sub("\S*\d\S*", "", sentance).strip()
    sentance = re.sub('[^A-Za-z]+', ' ', sentance)
    # https://gist.github.com/sebleier/554280
    sentance = ' '.join(e.lower() for e in sentance.split() if e.lower
() not in stopwords)
    preprocessed_reviews_100k.append(sentance.strip())
```

```
100%|████████████| 87773/87773 [00:34<00:00, 2548.09it/s]
```

In [0]:
```python
# Combining all the above stundents
from tqdm import tqdm
preprocessed_reviews_20k = []
# tqdm is for printing the status bar
for sentance in tqdm(final_20k['Text'].values):
    sentance = re.sub(r"http\S+", "", sentance)
    sentance = BeautifulSoup(sentance, 'lxml').get_text()
    sentance = decontracted(sentance)
    sentance = re.sub("\S*\d\S*", "", sentance).strip()
    sentance = re.sub('[^A-Za-z]+', ' ', sentance)
    # https://gist.github.com/sebleier/554280
    sentance = ' '.join(e.lower() for e in sentance.split() if e.lower
```

```
       () not in stopwords)
            preprocessed_reviews_20k.append(sentance.strip())
```

```
100%|████████| 19354/19354 [00:07<00:00, 2588.22it/s]
```

In [0]: `len(preprocessed_reviews_100k)`

Out[0]: 87773

In [0]: `len(preprocessed_reviews_20k)`

Out[0]: 19354

## [5] Assignment 7: SVM

1. **Apply SVM on these feature sets**

   - SET 1:Review text, preprocessed one converted into vectors using (BOW)
   - SET 2:Review text, preprocessed one converted into vectors using (TFIDF)
   - SET 3:Review text, preprocessed one converted into vectors using (AVG W2v)
   - SET 4:Review text, preprocessed one converted into vectors using (TFIDF W2v)

2. **Procedure**

   - You need to work with 2 versions of SVM
       - Linear kernel
       - RBF kernel
   - When you are working with linear kernel, use SGDClassifier' with hinge loss because it is computationally less expensive.
   - When you are working with 'SGDClassifier' with hinge loss and trying to find the AUC score, you would have to use CalibratedClassifierCV
   - Similarly, like kdtree of knn, when you are working with RBF kernel it's better to reduce the number of dimensions. You can put min_df = 10, max_features =

500 and consider a sample size of 40k points.

3. **Hyper paramter tuning (find best alpha in range [10^-4 to 10^4], and the best penalty among 'l1', 'l2')**

   - Find the best hyper parameter which will give the maximum [AUC](#) value
   - Find the best hyper paramter using k-fold cross validation or simple cross validation data
   - Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

4. **Feature importance**

   - When you are working on the linear kernel with BOW or TFIDF please print the top 10 best features for each of the positive and negative classes.

5. **Feature engineering**

   - To increase the performance of your model, you can also experiment with with feature engineering like :
     - Taking length of reviews as another feature.
     - Considering some features from review summary as well.

6. **Representation of results**

   - You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.
     Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.
     Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps.](#)
     

7. **[Conclusion](#)**

- [You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this prettytable library link](#)



**Note: Data Leakage**

1. There will be an issue of data-leakage if you vectorize the entire data and then split it into train/cv/test.
2. To avoid the issue of data-leakag, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method fit_transform() on you train data, and apply the method transform() on cv/test data.
4. For more details please go through this [link.](#)

# Applying SVM

## [5.1] Linear SVM

### [5.1.1] featurization for 100k dataset

```
In [0]:  #here preprocessed_review is my X and final['Score'] is my Y
         print(len(preprocessed_reviews_100k))
         print(len(final_100k['Score']))
         X=preprocessed_reviews_100k
         Y=final_100k['Score']
         #if both are of same lenght then proceed....
```

```
87773
87773
```

```
In [0]:  #here i am performing splittig operation as train test and cv...
         from sklearn.model_selection import train_test_split
```

```python
# X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=
0.33, shuffle=Flase)# this is for time series split
X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.3
3) # this is random splitting
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_
size=0.33) # this is random splitting
```

In [0]:
```python
#checking the types of test and train X,y
print(type(X_train))
print(type(X_test))
print(type(X_cv))
print(type(y_train))
print(type(y_test))
print(type(y_cv))
#now i have xtrain ,xtest,tcv and ytrain,ytest ,ycv....
```

```
<class 'list'>
<class 'list'>
<class 'list'>
<class 'pandas.core.series.Series'>
<class 'pandas.core.series.Series'>
<class 'pandas.core.series.Series'>
```

In [0]:
```python
#now you are ready with xtrain ,xtest xcv and ytrain ,ytest ,ycv
#now there is no problem to proceed with featurization
#first we will do Bow AND LATER OTHER...
```

**[5.1.1.1]BOW**

In [0]:
```python
#BoW
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer()
vectorizer.fit(X_train) # fitting on train data ,we cant perform fit on
 test or cv

# we use the fitted CountVectorizer to convert the text to vector
```

```
X_train_bow = vectorizer.transform(X_train)
X_cv_bow = vectorizer.transform(X_cv)
X_test_bow = vectorizer.transform(X_test)
print("After vectorizations")
print(X_train_bow.shape, y_train.shape)
print(X_cv_bow.shape, y_cv.shape)
print(X_test_bow.shape, y_test.shape)
print("="*100)
#you can also check X_train_bow is of sparse matrix type or not
#below is code for that
print(type(X_train_bow))
#displaying number of unique words in each of splitted dataset
print("the number of unique words in train: ", X_train_bow.get_shape()[
1])
print("the number of unique words in cv: ", X_cv_bow.get_shape()[1])
print("the number of unique words in test: ", X_test_bow.get_shape()[1
])
```

```
After vectorizations
(39400, 37791) (39400,)
(19407, 37791) (19407,)
(28966, 37791) (28966,)
=======================================================================
=============================
<class 'scipy.sparse.csr.csr_matrix'>
the number of unique words in train:  37791
the number of unique words in cv:  37791
the number of unique words in test:  37791
```

**[5.1.1.2]TF-IDF**

In [0]:
```
#below code for converting to tfidf
#i refered sample solution to write this code
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
tf_idf_vect.fit(X_train)
print("some sample features(unique words in the corpus)",tf_idf_vect.ge
t_feature_names()[0:10])
print('='*50)
```

```
X_train_tf_idf = tf_idf_vect.transform(X_train)
X_test_tf_idf = tf_idf_vect.transform(X_test)
X_cv_tf_idf = tf_idf_vect.transform(X_cv)
print("the type of count vectorizer ",type(X_train_tf_idf))
print("the shape of out text TFIDF vectorizer ",X_train_tf_idf.get_shap
e())
print("the number of unique words including both unigrams and bigrams "
, X_train_tf_idf.get_shape()[1])
```

```
some sample features(unique words in the corpus) ['ability', 'able', 'a
ble buy', 'able chew', 'able drink', 'able eat', 'able enjoy', 'able fi
nd', 'able get', 'able give']
==================================================
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer  (39400, 23526)
the number of unique words including both unigrams and bigrams  23526
```

**[5.1.1.3]AVG-W2V**

In [0]:
```
#in average w2v the output is of list form and here we write same code
 of all train ,test and cv
#this code is for train data:
# Train your own Word2Vec model using your own text corpus
i=0
list_of_sentance_train=[]
for sentance in X_train:
    list_of_sentance_train.append(sentance.split())


#training word2vect model
from gensim.models import Word2Vec
from gensim.models import KeyedVectors
# this line of code trains your w2v model on the give list of sentances
w2v_model=Word2Vec(list_of_sentance_train,min_count=5,size=50, workers=
4)
w2v_words = list(w2v_model.wv.vocab)
print("number of words that occured minimum 5 times ",len(w2v_words))
```

```python
print("sample words ", w2v_words[0:50])

#this is the actuall code to convert word2vect to avg w2v:
from tqdm import tqdm
import numpy as np
# average Word2Vec
# compute average word2vec for each review.
sent_vectors_train = []; # the avg-w2v for each sentence/review is stor
ed in this list
for sent in tqdm(list_of_sentance_train): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, yo
u might need to change this to 300 if you use google's w2v
    cnt_words =0; # num of words with a valid vector in the sentence/re
view
    for word in sent: # for each word in a review/sentence
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    sent_vectors_train.append(sent_vec)
sent_vectors_train = np.array(sent_vectors_train)
print(sent_vectors_train.shape)
print(sent_vectors_train[0])
```

```
0%|              | 133/39400 [00:00<00:29, 1324.79it/s]
```

```
number of words that occured minimum 5 times  12162
sample words  ['bounty', 'bars', 'consist', 'soft', 'almost', 'center',
'sweet', 'coconut', 'coated', 'outer', 'layer', 'milk', 'chocolate', 'e
asily', 'buy', 'counter', 'australia', 'not', 'case', 'mind', 'diet',
'rich', 'junk', 'food', 'yes', 'candy', 'bar', 'worth', 'going', 'troub
le', 'ordering', 'internet', 'note', 'amazon', 'purchase', 'boxes', 'va
rying', 'numbers', 'inside', 'crunch', 'think', 'package', 'offers', 'b
est', 'value', 'time', 'may', 'entirely', 'interests', 'health']
```

```
100%|████████| 39400/39400 [01:09<00:00, 567.13it/s]
```

```
(39400, 50)
[ 1.80089722e-01  2.05140825e-02 -8.05299217e-01  1.77457694e-01
```

```
       -5.97381066e-01 -7.02637962e-02  4.61338516e-02 -1.49315512e-01
        4.54844921e-01  4.03732935e-02 -1.45282104e-01 -5.74814765e-02
       -2.69850784e-02  3.43609117e-01 -2.42372032e-01 -5.77786456e-02
        4.73667569e-04  4.95357078e-01 -2.46512696e-01  2.69851262e-01
       -1.84349048e-01  2.98629502e-01  3.79519494e-01 -2.84680361e-03
        9.77460868e-01 -5.81388269e-01 -1.49062361e-01  4.39391380e-01
        1.54149523e-01  5.13092168e-02  1.85735057e-01  3.22970181e-01
       -4.40778378e-01  5.23877055e-01 -8.87539327e-01  4.35135524e-02
        5.43463729e-03 -4.45507933e-01 -2.70383744e-01  4.86540575e-01
        3.74177488e-01 -8.06842588e-02 -2.10020291e-01  3.47714017e-01
        5.39048511e-01  2.01946155e-01  1.08760666e-01 -6.57880929e-01
       -2.70598858e-01 -2.38995908e-01]
```

In [0]:
```python
#this code is for test data:
# Train your own Word2Vec model using your own text corpus
i=0
list_of_sentance_test=[]
for sentance in X_test:
    list_of_sentance_test.append(sentance.split())


#training word2vect model
from gensim.models import Word2Vec
from gensim.models import KeyedVectors
# this line of code trains your w2v model on the give list of sentances
#i made below two statement as comment to avoid data leakage problem
#w2v_model=Word2Vec(list_of_sentance_test,min_count=5,size=50, workers=4)
#w2v_words = list(w2v_model.wv.vocab)
print("number of words that occured minimum 5 times ",len(w2v_words))
print("sample words ", w2v_words[0:50])

#this is the actuall code to convert word2vect to avg w2v:
from tqdm import tqdm
import numpy as np
# average Word2Vec
# compute average word2vec for each review.
sent_vectors_test = []; # the avg-w2v for each sentence/review is stored in this list
```

```python
for sent in tqdm(list_of_sentance_test): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, yo
u might need to change this to 300 if you use google's w2v
    cnt_words =0; # num of words with a valid vector in the sentence/re
view
    for word in sent: # for each word in a review/sentence
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    sent_vectors_test.append(sent_vec)
sent_vectors_test = np.array(sent_vectors_test)
print(sent_vectors_test.shape)
print(sent_vectors_test[0])
```

```
  0%|              | 70/28966 [00:00<00:41, 694.87it/s]
```

```
number of words that occured minimum 5 times  12162
sample words  ['bounty', 'bars', 'consist', 'soft', 'almost', 'center',
'sweet', 'coconut', 'coated', 'outer', 'layer', 'milk', 'chocolate', 'e
asily', 'buy', 'counter', 'australia', 'not', 'case', 'mind', 'diet',
'rich', 'junk', 'food', 'yes', 'candy', 'bar', 'worth', 'going', 'troub
le', 'ordering', 'internet', 'note', 'amazon', 'purchase', 'boxes', 'va
rying', 'numbers', 'inside', 'crunch', 'think', 'package', 'offers', 'b
est', 'value', 'time', 'may', 'entirely', 'interests', 'health']
```

```
100%|██████████| 28966/28966 [00:50<00:00, 575.04it/s]
```

```
(28966, 50)
[ 0.28467143  0.06281107 -0.47944649  0.27945947  0.32130887  0.0962654
7
 -0.13775616  0.21194507  0.07661421 -0.53902863  0.17412668 -0.4878547
8
  0.12098106  0.42572142 -0.19986091  0.3243257   0.59554522  0.8785378
 -0.07091743  0.29742988 -0.29448097  0.26058289 -0.09931917 -0.4184489
  0.13339064 -0.24838281 -0.59720091  0.67425163  0.32461321  0.2587499
4
 -0.07073563  0.41527788 -0.63524983  0.192898   -0.39522407  0.0186022
9
```

```
  0.49156101 -0.2869714   0.05315959  0.50476489  0.72856728 -0.4407562
1
  0.32184575  0.55191699  0.81517271 -0.14343182  0.35068515 -0.1917308
8
  0.40096128  0.28092947]
```

In [0]:
```python
#this code is for cv data:
# Train your own Word2Vec model using your own text corpus
i=0
list_of_sentance_cv=[]
for sentance in X_cv:
    list_of_sentance_cv.append(sentance.split())


#training word2vect model
from gensim.models import Word2Vec
from gensim.models import KeyedVectors
# this line of code trains your w2v model on the give list of sentances
#w2v_model=Word2Vec(list_of_sentance_cv,min_count=5,size=50, workers=4)
#w2v_words = list(w2v_model.wv.vocab)
print("number of words that occured minimum 5 times ",len(w2v_words))
print("sample words ", w2v_words[0:50])

#this is the actuall code to convert word2vect to avg w2v:
from tqdm import tqdm
import numpy as np
# average Word2Vec
# compute average word2vec for each review.
sent_vectors_cv = []; # the avg-w2v for each sentence/review is stored
 in this list
for sent in tqdm(list_of_sentance_cv): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, yo
u might need to change this to 300 if you use google's w2v
    cnt_words =0; # num of words with a valid vector in the sentence/re
view
    for word in sent: # for each word in a review/sentence
        if word in w2v_words:
            vec = w2v_model.wv[word]
```

```
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    sent_vectors_cv.append(sent_vec)
sent_vectors_cv= np.array(sent_vectors_cv)
print(sent_vectors_cv.shape)
print(sent_vectors_cv[0])
```

```
  0%|              | 0/19407 [00:00<?, ?it/s]
```

```
number of words that occured minimum 5 times  12162
sample words  ['bounty', 'bars', 'consist', 'soft', 'almost', 'center',
'sweet', 'coconut', 'coated', 'outer', 'layer', 'milk', 'chocolate', 'e
asily', 'buy', 'counter', 'australia', 'not', 'case', 'mind', 'diet',
'rich', 'junk', 'food', 'yes', 'candy', 'bar', 'worth', 'going', 'troub
le', 'ordering', 'internet', 'note', 'amazon', 'purchase', 'boxes', 'va
rying', 'numbers', 'inside', 'crunch', 'think', 'package', 'offers', 'b
est', 'value', 'time', 'may', 'entirely', 'interests', 'health']
```

```
100%|███████████| 19407/19407 [00:33<00:00, 580.43it/s]
```

```
(19407, 50)
[-6.70679881e-01  1.29489416e-01 -4.00717816e-01  6.63609790e-01
 -1.19812183e-01  4.72399494e-01 -6.59335039e-02 -3.63703390e-01
  8.18900866e-03  1.55203498e-01 -2.04977452e-04  1.99777312e-01
  2.52543883e-01  7.29024756e-02 -2.71235509e-01 -3.22417541e-01
  2.37283467e-01 -2.14250649e-01 -2.49897607e-01  6.55582342e-01
 -4.25267784e-01  3.27082985e-01 -2.12470211e-01 -3.84439195e-01
  7.36364306e-01 -1.93072568e-01 -8.97169850e-02  3.36497250e-01
  2.64380837e-01  2.35288304e-01  7.66406389e-02  8.52452070e-01
 -1.69626742e-01  4.74021827e-01 -7.65199710e-01  2.81899138e-01
  6.11398019e-01 -5.55901841e-02 -3.73956042e-01  7.52873769e-01
  3.76777449e-02 -5.26155738e-01 -8.18421600e-01  9.41237990e-02
  6.80949122e-01 -1.10573763e-01  5.87993649e-01 -2.06107380e-01
 -4.23502602e-01  1.92218562e-01]
```

**[5.1.1.4]tfidf-w2v**

```
In [0]:  #this is for train data
         i=0
         list_of_sentance_train=[]
         for sentance in X_train:
             list_of_sentance_train.append(sentance.split())


         # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
         model = TfidfVectorizer()
         tf_idf_matrix = model.fit_transform(X_train)
         # we are converting a dictionary with word as a key, and the idf as a v
         alue
         dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))



         # TF-IDF weighted Word2Vec
         tfidf_feat = model.get_feature_names() # tfidf words/col-names
         # final_tf_idf is the sparse matrix with row= sentence, col=word and ce
         ll_val = tfidf

         tfidf_sent_vectors_train = []; # the tfidf-w2v for each sentence/review
          is stored in this list
         row=0;
         for sent in tqdm(list_of_sentance_train): # for each review/sentence
             sent_vec = np.zeros(50) # as word vectors are of zero length
             weight_sum =0; # num of words with a valid vector in the sentence/r
         eview
             for word in sent: # for each word in a review/sentence
                 if word in w2v_words and word in tfidf_feat:
                     vec = w2v_model.wv[word]
         #                tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
                     # to reduce the computation we are
                     # dictionary[word] = idf value of word in whole courpus
                     # sent.count(word) = tf valeus of word in this review
                     tf_idf = dictionary[word]*(sent.count(word)/len(sent))
                     sent_vec += (vec * tf_idf)
                     weight_sum += tf_idf
             if weight_sum != 0:
```

```python
        sent_vec /= weight_sum
    tfidf_sent_vectors_train.append(sent_vec)
    row += 1
tfidf_sent_vectors_train= np.array(sent_vectors_train)
print(tfidf_sent_vectors_train.shape)
print(tfidf_sent_vectors_train[0])
```

```
100%|████████████| 39400/39400 [14:14<00:00, 46.09it/s]
```

```
(39400, 50)
[ 1.80089722e-01  2.05140825e-02 -8.05299217e-01  1.77457694e-01
 -5.97381066e-01 -7.02637962e-02  4.61338516e-02 -1.49315512e-01
  4.54844921e-01  4.03732935e-02 -1.45282104e-01 -5.74814765e-02
 -2.69850784e-02  3.43609117e-01 -2.42372032e-01 -5.77786456e-02
  4.73667569e-04  4.95357078e-01 -2.46512696e-01  2.69851262e-01
 -1.84349048e-01  2.98629502e-01  3.79519494e-01 -2.84680361e-03
  9.77460868e-01 -5.81388269e-01 -1.49062361e-01  4.39391380e-01
  1.54149523e-01  5.13092168e-02  1.85735057e-01  3.22970181e-01
 -4.40778378e-01  5.23877055e-01 -8.87539327e-01  4.35135524e-02
  5.43463729e-03 -4.45507933e-01 -2.70383744e-01  4.86540575e-01
  3.74177488e-01 -8.06842588e-02 -2.10020291e-01  3.47714017e-01
  5.39048511e-01  2.01946155e-01  1.08760666e-01 -6.57880929e-01
 -2.70598858e-01 -2.38995908e-01]
```

In [0]:
```python
#this is for test data
i=0
list_of_sentance_test=[]
for sentance in X_test:
    list_of_sentance_test.append(sentance.split())


# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
#model = TfidfVectorizer()
tf_idf_matrix = model.transform(X_test)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
```

```python
# TF-IDF weighted Word2Vec
tfidf_feat = model.get_feature_names() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and ce
ll_val = tfidf

tfidf_sent_vectors_test = []; # the tfidf-w2v for each sentence/review
 is stored in this list
row=0;
for sent in tqdm(list_of_sentance_test): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum =0; # num of words with a valid vector in the sentence/r
eview
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
#            tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    tfidf_sent_vectors_test.append(sent_vec)
    row += 1
tfidf_sent_vectors_test= np.array(sent_vectors_test)
print(tfidf_sent_vectors_test.shape)
print(tfidf_sent_vectors_test[0])
```

```
100%|██████████| 28966/28966 [10:21<00:00, 52.24it/s]
```

```
(28966, 50)
[ 0.28467143  0.06281107 -0.47944649  0.27945947  0.32130887  0.0962654
7
 -0.13775616  0.21194507  0.07661421 -0.53902863  0.17412668 -0.4878547
8
  0.12098106  0.42572142 -0.19986091  0.3243257   0.59554522  0.8785378
 -0.07091743  0.29742988 -0.29448097  0.26058289 -0.09931917 -0.4184489
  0.13339064 -0.24838281 -0.59720091  0.67425163  0.32461321  0.2587499
```

```
4
 -0.07073563  0.41527788 -0.63524983  0.192898   -0.39522407  0.0186022
9
  0.49156101 -0.2869714   0.05315959  0.50476489  0.72856728 -0.4407562
1
  0.32184575  0.55191699  0.81517271 -0.14343182  0.35068515 -0.1917308
8
  0.40096128  0.28092947]
```

In [0]:
```python
#this is for cv data
i=0
list_of_sentance_cv=[]
for sentance in X_cv:
    list_of_sentance_cv.append(sentance.split())


# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
#model = TfidfVectorizer()
tf_idf_matrix = model.transform(X_cv)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))



# TF-IDF weighted Word2Vec
tfidf_feat = model.get_feature_names() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val = tfidf

tfidf_sent_vectors_cv = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(list_of_sentance_cv): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum =0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
```

```python
            if word in w2v_words and word in tfidf_feat:
                vec = w2v_model.wv[word]
#                tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
                # to reduce the computation we are
                # dictionary[word] = idf value of word in whole courpus
                # sent.count(word) = tf valeus of word in this review
                tf_idf = dictionary[word]*(sent.count(word)/len(sent))
                sent_vec += (vec * tf_idf)
                weight_sum += tf_idf
        if weight_sum != 0:
            sent_vec /= weight_sum
        tfidf_sent_vectors_cv.append(sent_vec)
        row += 1
tfidf_sent_vectors_cv= np.array(sent_vectors_cv)
print(tfidf_sent_vectors_cv.shape)
print(tfidf_sent_vectors_cv[0])
```

```
100%|██████████| 19407/19407 [06:57<00:00, 46.48it/s]
```

```
(19407, 50)
[-6.70679881e-01  1.29489416e-01 -4.00717816e-01  6.63609790e-01
 -1.19812183e-01  4.72399494e-01 -6.59335039e-02 -3.63703390e-01
  8.18900866e-03  1.55203498e-01 -2.04977452e-04  1.99777312e-01
  2.52543883e-01  7.29024756e-02 -2.71235509e-01 -3.22417541e-01
  2.37283467e-01 -2.14250649e-01 -2.49897607e-01  6.55582342e-01
 -4.25267784e-01  3.27082985e-01 -2.12470211e-01 -3.84439195e-01
  7.36364306e-01 -1.93072568e-01 -8.97169850e-02  3.36497250e-01
  2.64380837e-01  2.35288304e-01  7.66406389e-02  8.52452070e-01
 -1.69626742e-01  4.74021827e-01 -7.65199710e-01  2.81899138e-01
  6.11398019e-01 -5.55901841e-02 -3.73956042e-01  7.52873769e-01
  3.76777449e-02 -5.26155738e-01 -8.18421600e-01  9.41237990e-02
  6.80949122e-01 -1.10573763e-01  5.87993649e-01 -2.06107380e-01
 -4.23502602e-01  1.92218562e-01]
```

### [5.1.2] Applying Linear SVM on BOW, SET 1

```python
In [0]:  #finding best alhpa and with l1 regularization which maximises auc
         from sklearn.model_selection import train_test_split
```

```python
from sklearn.model_selection import learning_curve, GridSearchCV
from sklearn.linear_model import SGDClassifier
from sklearn.model_selection import  cross_val_score


clf=SGDClassifier(loss='hinge',class_weight='balanced')
param_grid={'alpha':[10**-4, 10**-2, 10**0, 10**2, 10**4],'penalty':['l
1']}
#Using GridSearchCV
model = GridSearchCV(clf,param_grid,scoring='roc_auc',cv=5,verbose=1,n_
jobs=1)
model.fit(X_train_bow, y_train)

print(model.best_estimator_)
print(model.score(X_test_bow, y_test))
```

Fitting 5 folds for each of 5 candidates, totalling 25 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

```
SGDClassifier(alpha=0.0001, average=False, class_weight='balanced',
        early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
        l1_ratio=0.15, learning_rate='optimal', loss='hinge', max_iter=N
one,
        n_iter=None, n_iter_no_change=5, n_jobs=None, penalty='l1',
        power_t=0.5, random_state=None, shuffle=True, tol=None,
        validation_fraction=0.1, verbose=0, warm_start=False)
0.9110651364831871
```

[Parallel(n_jobs=1)]: Done  25 out of  25 | elapsed:    4.4s finished

In [0]:
```python
#finding best alhpa and with l1 regularization which maximises auc
from sklearn.model_selection import train_test_split
from sklearn.model_selection import learning_curve, GridSearchCV
from sklearn.linear_model import SGDClassifier
from sklearn.model_selection import  cross_val_score


clf=SGDClassifier(loss='hinge',class_weight='balanced')
```

```python
param_grid={'alpha':[10**-4, 10**-2, 10**0, 10**2, 10**4],'penalty':['l
2']}
#Using GridSearchCV
model = GridSearchCV(clf,param_grid,scoring='roc_auc',cv=5,verbose=1,n_
jobs=1)
model.fit(X_train_bow, y_train)

print(model.best_estimator_)
print(model.score(X_test_bow, y_test))
```

Fitting 5 folds for each of 5 candidates, totalling 25 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent
workers.

```
SGDClassifier(alpha=0.01, average=False, class_weight='balanced',
        early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
        l1_ratio=0.15, learning_rate='optimal', loss='hinge', max_iter=N
one,
        n_iter=None, n_iter_no_change=5, n_jobs=None, penalty='l2',
        power_t=0.5, random_state=None, shuffle=True, tol=None,
        validation_fraction=0.1, verbose=0, warm_start=False)
0.9266699253729377
```

[Parallel(n_jobs=1)]: Done  25 out of  25 | elapsed:    3.1s finished

In [0]:
```python
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc
_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
from sklearn.linear_model import SGDClassifier
from sklearn.calibration import CalibratedClassifierCV

clf=SGDClassifier(alpha=0.01,penalty='l2',epsilon=0.1,loss='hinge',clas
s_weight='balanced')
clf.fit(X_train_bow, y_train)

# roc_auc_score(y_true, y_score) the 2nd parameter should be probabilit
y estimates of the positive class
# not the predicted outputs
```
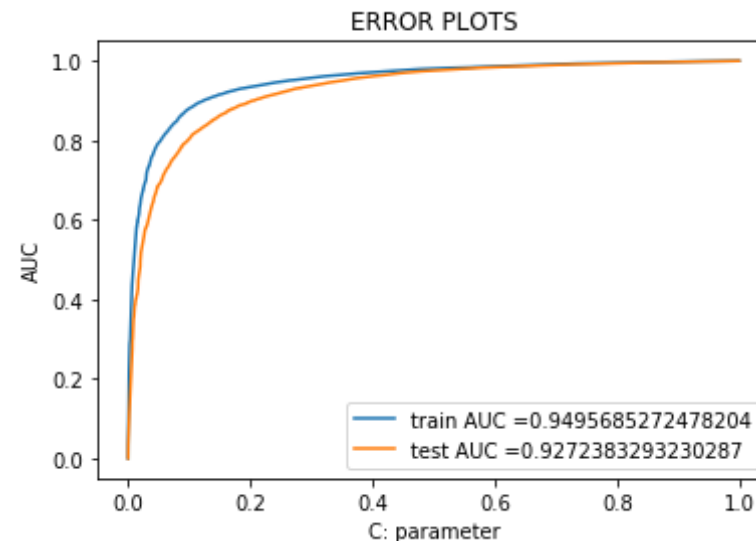
```python
#isotonic calibration
#https://scikit-learn.org/stable/auto_examples/calibration/plot_calibration.html
clf_isotonic = CalibratedClassifierCV(clf, cv=2, method='isotonic')
clf_isotonic.fit(X_train_bow, y_train)
#prob_pos_isotonic = clf_isotonic.predict_proba(X_test_bow)[:, 1]


train_fpr, train_tpr, thresholds = roc_curve(y_train, clf_isotonic.predict_proba(X_train_bow)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf_isotonic.predict_proba(X_test_bow)[:,1])

plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("C: parameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```

In [0]:
```python
#for seaborn confusion matrix :https://stackoverflow.com/questions/3557
2000/how-can-i-plot-a-confusion-matrix
#refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
fusion-matrix-with-labels
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
ax= plt.subplot()
arr1=confusion_matrix(y_train, clf.predict(X_train_bow))
df_1= pd.DataFrame(arr1, range(2),range(2))
plt.figure(figsize = (5,2))
sn.heatmap(df_1, annot=True,fmt="d",ax=ax)
ax.set_title('Confusion Matrix');
ax.set_xlabel('Actual Labels')
ax.set_ylabel('Predicted Labels')
ax.xaxis.set_ticklabels(['False', 'True']);
ax.yaxis.set_ticklabels(['True', 'False']);
```
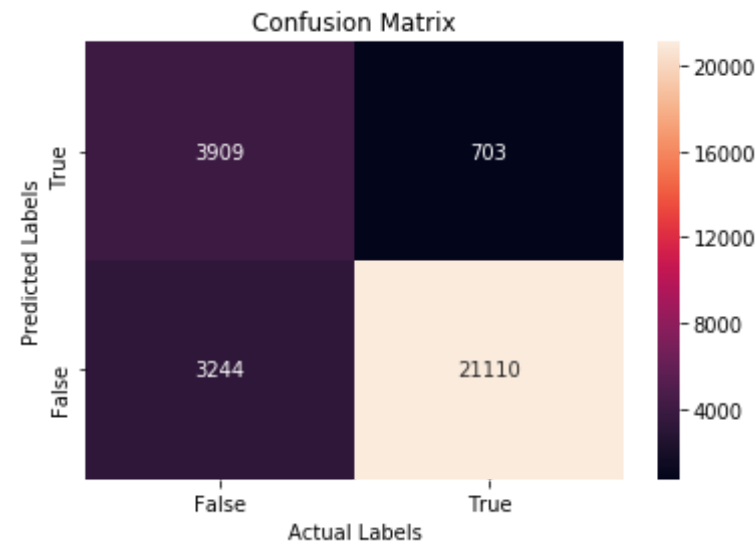
Train confusion matrix



<Figure size 360x144 with 0 Axes>

```
In [0]: #refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
        fusion-matrix-with-labels
        print("Test confusion matrix")
        ax= plt.subplot()
        arr2=confusion_matrix(y_test, clf.predict(X_test_bow))
        df_2= pd.DataFrame(arr2, range(2),range(2))
        plt.figure(figsize = (5,2))
        sn.heatmap(df_2, annot=True,fmt="d",ax=ax)
        ax.set_title('Confusion Matrix');
        ax.set_xlabel('Actual Labels')
        ax.set_ylabel('Predicted Labels')
        ax.xaxis.set_ticklabels(['False', 'True']);
        ax.yaxis.set_ticklabels(['True', 'False']);
```

Test confusion matrix



```
<Figure size 360x144 with 0 Axes>
```

```
In [0]: #https://medium.com/@aneesha/visualising-top-features-in-linear-svm-wi
        th-scikit-learn-and-matplotlib-3454ab18a14d
        coef = clf.coef_.ravel()
```

```
top_positive_coefficients= np.argsort(coef)[-10:]
top_negative_coefficients = np.argsort(coef)[:10]
print('top positive:')
print(np.take(vectorizer.get_feature_names(),top_positive_coefficients
))
print('top negetive:')
print(np.take(vectorizer.get_feature_names(),top_negative_coefficients
))
```

```
top positive:
['love' 'highly' 'nice' 'wonderful' 'excellent' 'perfect' 'loves' 'bes
t'
 'delicious' 'great']
top negetive:
['disappointed' 'worst' 'unfortunately' 'money' 'terrible' 'awful'
 'horrible' 'disappointing' 'thought' 'bad']
```

### [5.1.3] Applying Linear SVM on TFIDF, SET 2

In [0]:
```
#finding best alhpa and with l1 regularization which maximises auc
from sklearn.model_selection import train_test_split
from sklearn.model_selection import learning_curve, GridSearchCV
from sklearn.linear_model import SGDClassifier
from sklearn.model_selection import  cross_val_score


clf=SGDClassifier(loss='hinge',class_weight='balanced')
param_grid={'alpha':[10**-4, 10**-2, 10**0, 10**2, 10**4],'penalty':['l
1']}
#Using GridSearchCV
model = GridSearchCV(clf,param_grid,scoring='roc_auc',cv=5,verbose=1,n_
jobs=1)
model.fit(X_train_tf_idf, y_train)

print(model.best_estimator_)
print(model.score(X_test_tf_idf, y_test))
```

```
Fitting 5 folds for each of 5 candidates, totalling 25 fits
```

```
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent
workers.
[Parallel(n_jobs=1)]: Done  25 out of  25 | elapsed:     4.8s finished
```

```
SGDClassifier(alpha=0.0001, average=False, class_weight='balanced',
        early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
        l1_ratio=0.15, learning_rate='optimal', loss='hinge', max_iter=N
one,
        n_iter=None, n_iter_no_change=5, n_jobs=None, penalty='l1',
        power_t=0.5, random_state=None, shuffle=True, tol=None,
        validation_fraction=0.1, verbose=0, warm_start=False)
0.9442385428545605
```

In [0]:
```python
#finding best alhpa and with l2 regularization which maximises auc
from sklearn.model_selection import train_test_split
from sklearn.model_selection import learning_curve, GridSearchCV
from sklearn.linear_model import SGDClassifier
from sklearn.model_selection import  cross_val_score


clf=SGDClassifier(loss='hinge',class_weight='balanced')
param_grid={'alpha':[10**-4, 10**-2, 10**0, 10**2, 10**4],'penalty':['l
2']}
#Using GridSearchCV
model = GridSearchCV(clf,param_grid,scoring='roc_auc',cv=5,verbose=1,n_
jobs=1)
model.fit(X_train_tf_idf, y_train)

print(model.best_estimator_)
print(model.score(X_test_tf_idf, y_test))
```

```
Fitting 5 folds for each of 5 candidates, totalling 25 fits
```

```
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent
workers.
```

```
SGDClassifier(alpha=0.0001, average=False, class_weight='balanced',
        early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
        l1_ratio=0.15, learning_rate='optimal', loss='hinge', max_iter=N
one,
```

```
                       n_iter=None, n_iter_no_change=5, n_jobs=None, penalty='l2',
                       power_t=0.5, random_state=None, shuffle=True, tol=None,
                       validation_fraction=0.1, verbose=0, warm_start=False)
          0.9579138423417927

          [Parallel(n_jobs=1)]: Done  25 out of  25 | elapsed:    3.9s finished
```

In [0]:
```python
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc
_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
from sklearn.linear_model import SGDClassifier
from sklearn.calibration import CalibratedClassifierCV

clf=SGDClassifier(alpha=0.0001,penalty='l2',epsilon=0.1,loss='hinge',cl
ass_weight='balanced')
clf.fit(X_train_tf_idf, y_train)

# roc_auc_score(y_true, y_score) the 2nd parameter should be probabilit
y estimates of the positive class
# not the predicted outputs

#isotonic calibration
#https://scikit-learn.org/stable/auto_examples/calibration/plot_calibra
tion.html
clf_isotonic = CalibratedClassifierCV(clf, cv=2, method='isotonic')
clf_isotonic.fit(X_train_tf_idf, y_train)


train_fpr, train_tpr, thresholds = roc_curve(y_train, clf_isotonic.pred
ict_proba(X_train_tf_idf)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf_isotonic.predict
_proba(X_test_tf_idf)[:,1])

plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, t
rain_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_
tpr)))
plt.legend()
plt.xlabel("C: parameter")
```

```
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```



In [0]:
```
#for seaborn confusion matrix :https://stackoverflow.com/questions/3557
2000/how-can-i-plot-a-confusion-matrix
#refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
fusion-matrix-with-labels
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
ax= plt.subplot()
arr1=confusion_matrix(y_train, clf.predict(X_train_tf_idf))
df_1= pd.DataFrame(arr1, range(2),range(2))
plt.figure(figsize = (5,2))
sn.heatmap(df_1, annot=True,fmt="d",ax=ax)
ax.set_title('Confusion Matrix');
ax.set_xlabel('Actual Labels')
ax.set_ylabel('Predicted Labels')
```
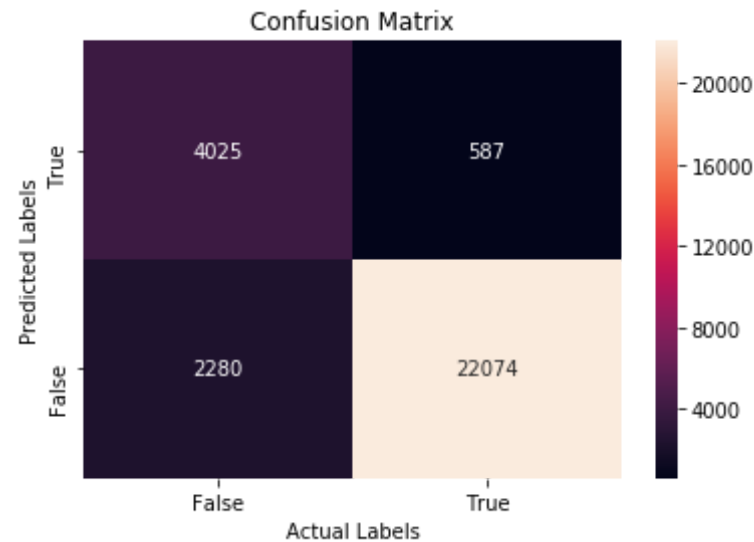
```
ax.xaxis.set_ticklabels(['False', 'True']);
ax.yaxis.set_ticklabels(['True', 'False']);
```

Train confusion matrix



<Figure size 360x144 with 0 Axes>

In [0]:
```
#refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
fusion-matrix-with-labels
print("Test confusion matrix")
ax= plt.subplot()
arr2=confusion_matrix(y_test, clf.predict(X_test_tf_idf))
df_2= pd.DataFrame(arr2, range(2),range(2))
plt.figure(figsize = (5,2))
sn.heatmap(df_2, annot=True,fmt="d",ax=ax)
ax.set_title('Confusion Matrix');
ax.set_xlabel('Actual Labels')
ax.set_ylabel('Predicted Labels')
ax.xaxis.set_ticklabels(['False', 'True']);
ax.yaxis.set_ticklabels(['True', 'False']);
```

Test confusion matrix

```
<Figure size 360x144 with 0 Axes>
```

**[5.1.4] Applying Linear SVM on AVG W2V, SET 3**

```python
In [0]:  #finding best alhpa and with l1 regularization which maximises auc
         from sklearn.model_selection import train_test_split
         from sklearn.model_selection import learning_curve, GridSearchCV
         from sklearn.linear_model import SGDClassifier
         from sklearn.model_selection import  cross_val_score


         clf=SGDClassifier(loss='hinge',class_weight='balanced')
         param_grid={'alpha':[10**-4, 10**-2, 10**0, 10**2, 10**4],'penalty':['l
         1']}
         #Using GridSearchCV
         model = GridSearchCV(clf,param_grid,scoring='roc_auc',cv=5,verbose=1,n_
         jobs=1)
         model.fit(sent_vectors_train, y_train)

         print(model.best_estimator_)
         print(model.score(sent_vectors_test, y_test))
```

```
Fitting 5 folds for each of 5 candidates, totalling 25 fits
```

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

```
SGDClassifier(alpha=0.01, average=False, class_weight='balanced',
        early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
        l1_ratio=0.15, learning_rate='optimal', loss='hinge', max_iter=None,
        n_iter=None, n_iter_no_change=5, n_jobs=None, penalty='l1',
        power_t=0.5, random_state=None, shuffle=True, tol=None,
        validation_fraction=0.1, verbose=0, warm_start=False)
0.8869022238903038
```

[Parallel(n_jobs=1)]: Done  25 out of  25 | elapsed:    4.7s finished

In [0]:
```python
#finding best alhpa and with l2 regularization which maximises auc
from sklearn.model_selection import train_test_split
from sklearn.model_selection import learning_curve, GridSearchCV
from sklearn.linear_model import SGDClassifier
from sklearn.model_selection import  cross_val_score


clf=SGDClassifier(loss='hinge',class_weight='balanced')
param_grid={'alpha':[10**-4, 10**-2, 10**0, 10**2, 10**4],'penalty':['l2']}
#Using GridSearchCV
model = GridSearchCV(clf,param_grid,scoring='roc_auc',cv=5,verbose=1,n_jobs=1)
model.fit(sent_vectors_train, y_train)

print(model.best_estimator_)
print(model.score(sent_vectors_test, y_test))
```

```
Fitting 5 folds for each of 5 candidates, totalling 25 fits
```

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

```
SGDClassifier(alpha=0.01, average=False, class_weight='balanced',
```

```
        early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
        l1_ratio=0.15, learning_rate='optimal', loss='hinge', max_iter=N
one,
        n_iter=None, n_iter_no_change=5, n_jobs=None, penalty='l2',
        power_t=0.5, random_state=None, shuffle=True, tol=None,
        validation_fraction=0.1, verbose=0, warm_start=False)
0.8925402522606529
```

```
[Parallel(n_jobs=1)]: Done  25 out of  25 | elapsed:    2.8s finished
```

In [0]:
```python
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc
_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
from sklearn.linear_model import SGDClassifier
from sklearn.calibration import CalibratedClassifierCV

clf=SGDClassifier(alpha=0.0001,penalty='l2',epsilon=0.1,loss='hinge',cl
ass_weight='balanced')
clf.fit(sent_vectors_train, y_train)

# roc_auc_score(y_true, y_score) the 2nd parameter should be probabilit
y estimates of the positive class
# not the predicted outputs

#isotonic calibration
#https://scikit-learn.org/stable/auto_examples/calibration/plot_calibra
tion.html
clf_isotonic = CalibratedClassifierCV(clf, cv=2, method='isotonic')
clf_isotonic.fit(sent_vectors_train, y_train)


train_fpr, train_tpr, thresholds = roc_curve(y_train, clf_isotonic.pred
ict_proba(sent_vectors_train)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf_isotonic.predict
_proba(sent_vectors_test)[:,1])

plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, t
rain_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_
```
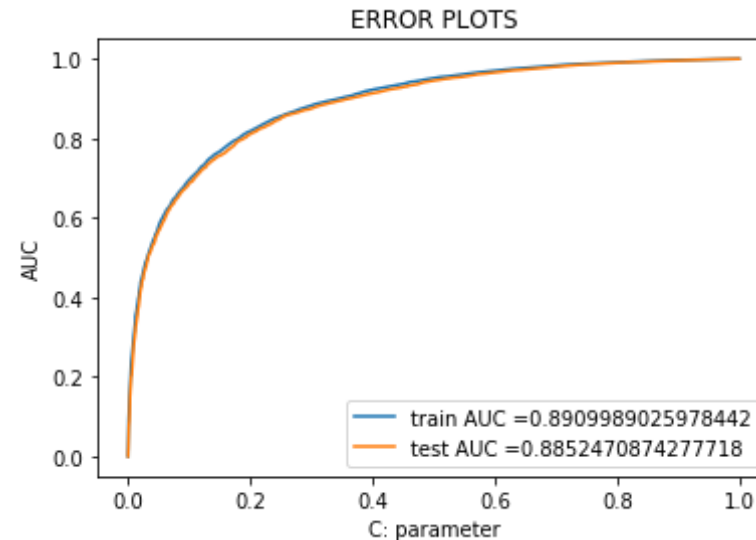
```
        tpr)))
plt.legend()
plt.xlabel("C: parameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```
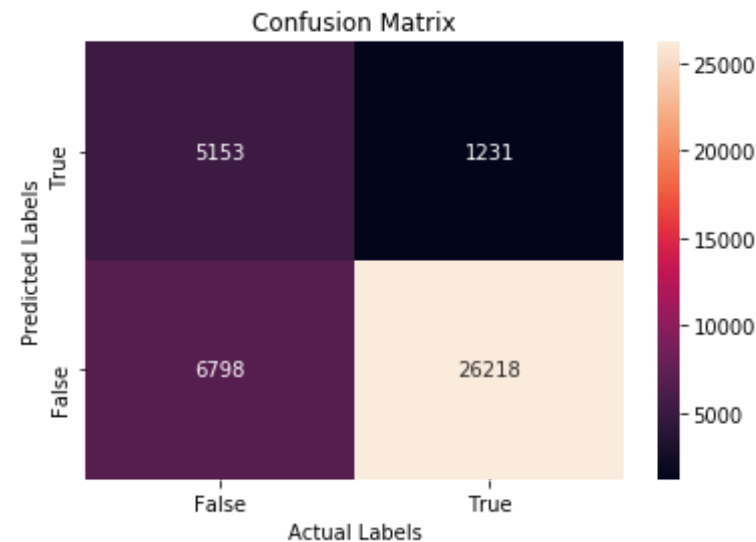
ERROR PLOTS



In [0]:
```
#for seaborn confusion matrix :https://stackoverflow.com/questions/3557
2000/how-can-i-plot-a-confusion-matrix
#refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
fusion-matrix-with-labels
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
ax= plt.subplot()
arr1=confusion_matrix(y_train, clf.predict(sent_vectors_train))
df_1= pd.DataFrame(arr1, range(2),range(2))
plt.figure(figsize = (5,2))
sn.heatmap(df_1, annot=True,fmt="d",ax=ax)
ax.set_title('Confusion Matrix');
```

```
ax.set_xlabel('Actual Labels')
ax.set_ylabel('Predicted Labels')
ax.xaxis.set_ticklabels(['False', 'True']);
ax.yaxis.set_ticklabels(['True', 'False']);
```
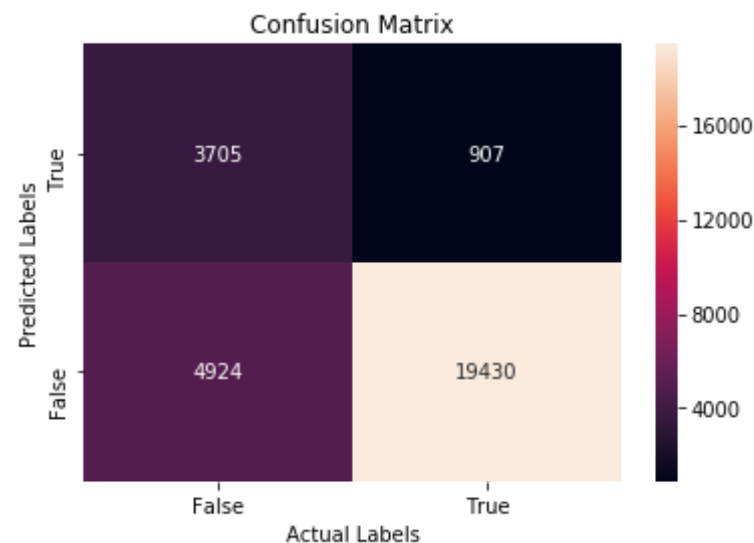
Train confusion matrix



```
<Figure size 360x144 with 0 Axes>
```

In [0]:
```
#refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
fusion-matrix-with-labels
print("Test confusion matrix")
ax= plt.subplot()
arr2=confusion_matrix(y_test, clf.predict(sent_vectors_test))
df_2= pd.DataFrame(arr2, range(2),range(2))
plt.figure(figsize = (5,2))
sn.heatmap(df_2, annot=True,fmt="d",ax=ax)
ax.set_title('Confusion Matrix');
ax.set_xlabel('Actual Labels')
ax.set_ylabel('Predicted Labels')
ax.xaxis.set_ticklabels(['False', 'True']);
ax.yaxis.set_ticklabels(['True', 'False']);
```

Test confusion matrix


Confusion Matrix

<Figure size 360x144 with 0 Axes>

### [5.1.5] Applying Linear SVM on TFIDF W2V, <span style="color:red">SET 4</span>

```python
In [0]:  #finding best alhpa and with l1 regularization which maximises auc
         from sklearn.model_selection import train_test_split
         from sklearn.model_selection import learning_curve, GridSearchCV
         from sklearn.linear_model import SGDClassifier
         from sklearn.model_selection import  cross_val_score


         clf=SGDClassifier(loss='hinge',class_weight='balanced')
         param_grid={'alpha':[10**-4, 10**-2, 10**0, 10**2, 10**4],'penalty':['l
         1']}
         #Using GridSearchCV
         model = GridSearchCV(clf,param_grid,scoring='roc_auc',cv=5,verbose=1,n_
         jobs=1)
         model.fit(tfidf_sent_vectors_train, y_train)
```

```
print(model.best_estimator_)
print(model.score(tfidf_sent_vectors_test, y_test))
```

Fitting 5 folds for each of 5 candidates, totalling 25 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent
workers.

```
SGDClassifier(alpha=0.01, average=False, class_weight='balanced',
        early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
        l1_ratio=0.15, learning_rate='optimal', loss='hinge', max_iter=N
one,
        n_iter=None, n_iter_no_change=5, n_jobs=None, penalty='l1',
        power_t=0.5, random_state=None, shuffle=True, tol=None,
        validation_fraction=0.1, verbose=0, warm_start=False)
0.8867568810678513
```

[Parallel(n_jobs=1)]: Done  25 out of  25 | elapsed:    4.7s finished

In [0]:
```
#finding best alhpa and with l1 regularization which maximises auc
from sklearn.model_selection import train_test_split
from sklearn.model_selection import learning_curve, GridSearchCV
from sklearn.linear_model import SGDClassifier
from sklearn.model_selection import  cross_val_score


clf=SGDClassifier(loss='hinge',class_weight='balanced')
param_grid={'alpha':[10**-4, 10**-2, 10**0, 10**2, 10**4],'penalty':['l
2']}
#Using GridSearchCV
model = GridSearchCV(clf,param_grid,scoring='roc_auc',cv=5,verbose=1,n_
jobs=1)
model.fit(tfidf_sent_vectors_train, y_train)

print(model.best_estimator_)
print(model.score(tfidf_sent_vectors_test, y_test))
```

Fitting 5 folds for each of 5 candidates, totalling 25 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent
workers.

```
SGDClassifier(alpha=0.01, average=False, class_weight='balanced',
        early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
        l1_ratio=0.15, learning_rate='optimal', loss='hinge', max_iter=N
one,
        n_iter=None, n_iter_no_change=5, n_jobs=None, penalty='l2',
        power_t=0.5, random_state=None, shuffle=True, tol=None,
        validation_fraction=0.1, verbose=0, warm_start=False)
0.8926237898841183
```

```
[Parallel(n_jobs=1)]: Done  25 out of  25 | elapsed:    3.2s finished
```

In [0]:
```python
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc
_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
from sklearn.linear_model import SGDClassifier
from sklearn.calibration import CalibratedClassifierCV

clf=SGDClassifier(alpha=0.01,penalty='l2',epsilon=0.1,loss='hinge',clas
s_weight='balanced')
clf.fit(sent_vectors_train, y_train)

# roc_auc_score(y_true, y_score) the 2nd parameter should be probabilit
y estimates of the positive class
# not the predicted outputs

#isotonic calibration
#https://scikit-learn.org/stable/auto_examples/calibration/plot_calibra
tion.html
clf_isotonic = CalibratedClassifierCV(clf, cv=2, method='isotonic')
clf_isotonic.fit(tfidf_sent_vectors_train, y_train)


train_fpr, train_tpr, thresholds = roc_curve(y_train, clf_isotonic.pred
ict_proba(tfidf_sent_vectors_train)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf_isotonic.predict
_proba(tfidf_sent_vectors_test)[:,1])

plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, t
rain_tpr)))
```
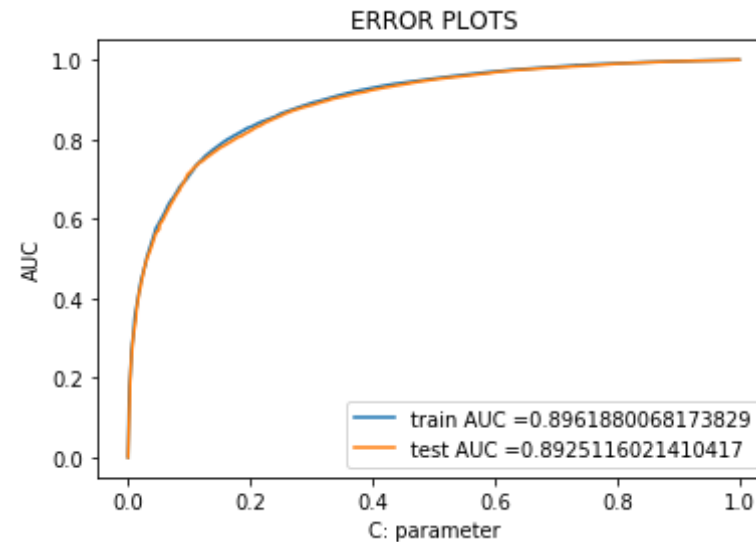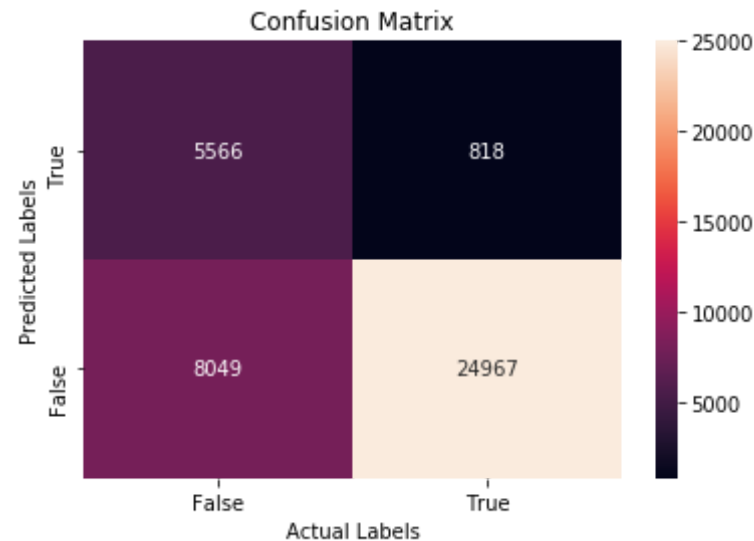
```
plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_
tpr)))
plt.legend()
plt.xlabel("C: parameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```

ERROR PLOTS



In [0]:
```
#for seaborn confusion matrix :https://stackoverflow.com/questions/3557
2000/how-can-i-plot-a-confusion-matrix
#refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
fusion-matrix-with-labels
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
ax= plt.subplot()
arr1=confusion_matrix(y_train, clf.predict(tfidf_sent_vectors_train))
df_1= pd.DataFrame(arr1, range(2),range(2))
plt.figure(figsize = (5,2))
sn.heatmap(df_1, annot=True,fmt="d",ax=ax)
```

```
ax.set_title('Confusion Matrix');
ax.set_xlabel('Actual Labels')
ax.set_ylabel('Predicted Labels')
ax.xaxis.set_ticklabels(['False', 'True']);
ax.yaxis.set_ticklabels(['True', 'False']);
```
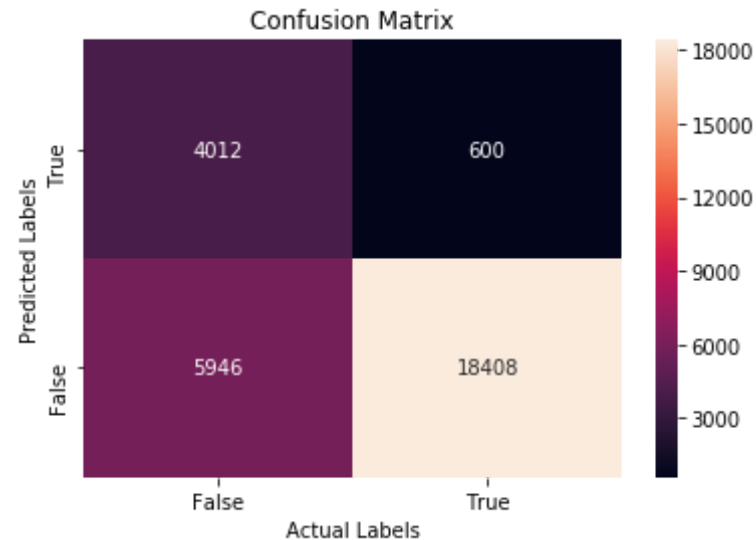
Train confusion matrix



```
<Figure size 360x144 with 0 Axes>
```

In [0]:
```
#refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
fusion-matrix-with-labels
print("Test confusion matrix")
ax= plt.subplot()
arr2=confusion_matrix(y_test, clf.predict(tfidf_sent_vectors_test))
df_2= pd.DataFrame(arr2, range(2),range(2))
plt.figure(figsize = (5,2))
sn.heatmap(df_2, annot=True,fmt="d",ax=ax)
ax.set_title('Confusion Matrix');
ax.set_xlabel('Actual Labels')
ax.set_ylabel('Predicted Labels')
ax.xaxis.set_ticklabels(['False', 'True']);
ax.yaxis.set_ticklabels(['True', 'False']);
```

Test confusion matrix



Confusion Matrix

<Figure size 360x144 with 0 Axes>

## [5.2] RBF SVM

### [5.2.1]featurization of 20k dataset

```
In [0]:  #here preprocessed_review is my X and final['Score'] is my Y
         print(len(preprocessed_reviews_20k))
         print(len(final_20k['Score']))
         X=preprocessed_reviews_20k
         Y=final_20k['Score']
         #if both are of same lenght then proceed....
```

```
19354
19354
```

```
In [0]:  #here i am performing splittig operation as train test and cv...
```

```python
from sklearn.model_selection import train_test_split

# X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=
0.33, shuffle=Flase)# this is for time series split
X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.3
3) # this is random splitting
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_
size=0.33) # this is random splitting
```

**[5.2.1.1]BOW**

In [0]:
```python
#BoW
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(min_df = 10, max_features = 500)
vectorizer.fit(X_train) # fitting on train data ,we cant perform fit on
 test or cv

# we use the fitted CountVectorizer to convert the text to vector
X_train_bow = vectorizer.transform(X_train)
X_cv_bow = vectorizer.transform(X_cv)
X_test_bow = vectorizer.transform(X_test)
print("After vectorizations")
print(X_train_bow.shape, y_train.shape)
print(X_cv_bow.shape, y_cv.shape)
print(X_test_bow.shape, y_test.shape)
print("="*100)
#you can also check X_train_bow is of sparse matrix type or not
#below is code for that
print(type(X_train_bow))
#displaying number of unique words in each of splitted dataset
print("the number of unique words in train: ", X_train_bow.get_shape()[
1])
print("the number of unique words in cv: ", X_cv_bow.get_shape()[1])
print("the number of unique words in test: ", X_test_bow.get_shape()[1
])
```

```
After vectorizations
(8687, 500) (8687,)
```

```
(4280, 500) (4280,)
(6387, 500) (6387,)
================================================================================
==============================
<class 'scipy.sparse.csr.csr_matrix'>
the number of unique words in train:  500
the number of unique words in cv:  500
the number of unique words in test:  500
```

**[5.2.1.2]TFIDF**

In [0]:
```python
#below code for converting to tfidf
#i refered sample solution to write this code
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2),min_df = 10, max_featur
es = 500)
tf_idf_vect.fit(X_train)
print("some sample features(unique words in the corpus)",tf_idf_vect.ge
t_feature_names()[0:10])
print('='*50)

X_train_tf_idf = tf_idf_vect.transform(X_train)
X_test_tf_idf = tf_idf_vect.transform(X_test)
X_cv_tf_idf = tf_idf_vect.transform(X_cv)
print("the type of count vectorizer ",type(X_train_tf_idf))
print("the shape of out text TFIDF vectorizer ",X_train_tf_idf.get_shap
e())
print("the number of unique words including both unigrams and bigrams "
, X_train_tf_idf.get_shape()[1])
```

```
some sample features(unique words in the corpus) ['able', 'absolutely',
'actually', 'add', 'added', 'aftertaste', 'ago', 'almonds', 'almost',
'along']
==================================================
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer  (8687, 500)
the number of unique words including both unigrams and bigrams  500
```

**[5.2.1.3]avg w2v**

In [0]:
```python
#in average w2v the output is of list form and here we write same code
 of all train ,test and cv
#this code is for train data:
# Train your own Word2Vec model using your own text corpus
i=0
list_of_sentance_train=[]
for sentence in X_train:
    list_of_sentance_train.append(sentance.split())


#training word2vect model
from gensim.models import Word2Vec
from gensim.models import KeyedVectors
# this line of code trains your w2v model on the give list of sentances
w2v_model=Word2Vec(list_of_sentance_train,min_count=5,size=50, workers=
4)
w2v_words = list(w2v_model.wv.vocab)
print("number of words that occured minimum 5 times ",len(w2v_words))
print("sample words ", w2v_words[0:50])

#this is the actuall code to convert word2vect to avg w2v:
from tqdm import tqdm
import numpy as np
# average Word2Vec
# compute average word2vec for each review.
sent_vectors_train = []; # the avg-w2v for each sentence/review is stor
ed in this list
for sent in tqdm(list_of_sentance_train): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, yo
u might need to change this to 300 if you use google's w2v
    cnt_words =0; # num of words with a valid vector in the sentence/re
view
    for word in sent: # for each word in a review/sentence
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
```

```
    if cnt_words != 0:
        sent_vec /= cnt_words
    sent_vectors_train.append(sent_vec)
sent_vectors_train = np.array(sent_vectors_train)
print(sent_vectors_train.shape)
print(sent_vectors_train[0])
```

```
  2%||              | 142/8687 [00:00<00:06, 1412.49it/s]
```

```
number of words that occured minimum 5 times  5495
sample words  ['love', 'pineapple', 'flavor', 'makes', 'feel', 'like',
'vacation', 'drinking', 'pina', 'probably', 'could', 'add', 'little',
'rum', 'really', 'tastes', 'better', 'coconut', 'waters', 'tried', 'sta
rted', 'son', 'baby', 'food', 'month', 'birthday', 'almost', 'months',
'introducing', 'avocado', 'quinoa', 'cereal', 'squash', 'huge', 'hit',
'sweet', 'potatoes', 'peas', 'apples', 'organic', 'taste', 'fresh', 'mu
ch', 'prefer', 'make', 'give', 'jarred', 'often', 'use', 'different']
```

```
100%|████████████| 8687/8687 [00:09<00:00, 874.39it/s]
```

```
(8687, 50)
[ 0.21411499 -0.36776797 -0.75282662  0.33994114  0.40363468 -0.4364555
3
 -0.21620836 -0.39892652  0.31709211  0.3927634   0.31608421  0.2378958
1
  0.15208639  0.5907577  -0.59054357 -0.43573927  0.16540512  0.1558099
8
 -0.10617807  0.7075104  -0.6076003   0.54247047  0.00242196 -0.7233496
5
  0.53540209 -0.20284717 -0.0522027   0.04519267 -0.40957193  0.5115167
9
  0.29258774  0.76476256 -0.44577942  0.44159299 -0.44819606  0.3085005
3
  0.31808708  0.0550734  -0.35594978 -0.03504065 -0.07206956 -0.2751188
3
 -0.05680525  0.21498388  0.59119366  0.0805137   0.4885871  -0.5601444
2
 -0.00090877 -0.41330399]
```

In [0]: `#this code is for test data:`

```python
# Train your own Word2Vec model using your own text corpus
i=0
list_of_sentance_test=[]
for sentance in X_test:
    list_of_sentance_test.append(sentance.split())


#training word2vect model
from gensim.models import Word2Vec
from gensim.models import KeyedVectors
# this line of code trains your w2v model on the give list of sentances
#i made below two statement as comment to avoid data leakage problem
#w2v_model=Word2Vec(list_of_sentance_test,min_count=5,size=50, workers=4)
#w2v_words = list(w2v_model.wv.vocab)
print("number of words that occured minimum 5 times ",len(w2v_words))
print("sample words ", w2v_words[0:50])

#this is the actuall code to convert word2vect to avg w2v:
from tqdm import tqdm
import numpy as np
# average Word2Vec
# compute average word2vec for each review.
sent_vectors_test = []; # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(list_of_sentance_test): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to change this to 300 if you use google's w2v
    cnt_words =0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    sent_vectors_test.append(sent_vec)
sent_vectors_test = np.array(sent_vectors_test)
```

```
print(sent_vectors_test.shape)
print(sent_vectors_test[0])
```

  2%|                 | 103/6387 [00:00<00:06, 1019.08it/s]

number of words that occured minimum 5 times  5495
sample words  ['love', 'pineapple', 'flavor', 'makes', 'feel', 'like',
'vacation', 'drinking', 'pina', 'probably', 'could', 'add', 'little',
'rum', 'really', 'tastes', 'better', 'coconut', 'waters', 'tried', 'sta
rted', 'son', 'baby', 'food', 'month', 'birthday', 'almost', 'months',
'introducing', 'avocado', 'quinoa', 'cereal', 'squash', 'huge', 'hit',
'sweet', 'potatoes', 'peas', 'apples', 'organic', 'taste', 'fresh', 'mu
ch', 'prefer', 'make', 'give', 'jarred', 'often', 'use', 'different']

100%|███████████| 6387/6387 [00:07<00:00, 846.75it/s]

```
(6387, 50)
[ 0.1020971  -0.25778366 -0.65248146  0.4684123   0.31574818 -0.4129766
  5
 -0.20247924 -0.40664902  0.19689101  0.31480956  0.27975016  0.3556833
  4
  0.26890199  0.6044926  -0.56634632 -0.39209652  0.14483991  0.0565501
  4
  0.0810383   0.62992669 -0.66624729  0.55472585 -0.19655463 -0.7660524
  2
  0.66199181 -0.15676928 -0.17650486 -0.0201654  -0.14755801  0.5113350
  9
  0.20591084  0.8363934  -0.4064      0.41434614 -0.48395601  0.2026638
  8
  0.45174919  0.05811068 -0.35995091  0.11331295 -0.02787627 -0.4291904
  1
 -0.06195572  0.35714837  0.65272359 -0.02911503  0.38616773 -0.5571473
  0.07343648 -0.32804483]
```

In [0]:
```
#this code is for cv data:
# Train your own Word2Vec model using your own text corpus
i=0
list_of_sentance_cv=[]
for sentance in X_cv:
    list_of_sentance_cv.append(sentance.split())
```

```python
#training word2vect model
from gensim.models import Word2Vec
from gensim.models import KeyedVectors
# this line of code trains your w2v model on the give list of sentances
#w2v_model=Word2Vec(list_of_sentance_cv,min_count=5,size=50, workers=4)
#w2v_words = list(w2v_model.wv.vocab)
print("number of words that occured minimum 5 times ",len(w2v_words))
print("sample words ", w2v_words[0:50])

#this is the actuall code to convert word2vect to avg w2v:
from tqdm import tqdm
import numpy as np
# average Word2Vec
# compute average word2vec for each review.
sent_vectors_cv = []; # the avg-w2v for each sentence/review is stored
 in this list
for sent in tqdm(list_of_sentance_cv): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, yo
u might need to change this to 300 if you use google's w2v
    cnt_words =0; # num of words with a valid vector in the sentence/re
view
    for word in sent: # for each word in a review/sentence
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    sent_vectors_cv.append(sent_vec)
sent_vectors_cv= np.array(sent_vectors_cv)
print(sent_vectors_cv.shape)
print(sent_vectors_cv[0])
```

```
  2%||              | 91/4280 [00:00<00:04, 878.55it/s]
```

```
number of words that occured minimum 5 times  5495
sample words  ['love', 'pineapple', 'flavor', 'makes', 'feel', 'like',
'vacation', 'drinking', 'pina', 'probably', 'could', 'add', 'little',
```

```
'rum', 'really', 'tastes', 'better', 'coconut', 'waters', 'tried', 'sta
rted', 'son', 'baby', 'food', 'month', 'birthday', 'almost', 'months',
'introducing', 'avocado', 'quinoa', 'cereal', 'squash', 'huge', 'hit',
'sweet', 'potatoes', 'peas', 'apples', 'organic', 'taste', 'fresh', 'mu
ch', 'prefer', 'make', 'give', 'jarred', 'often', 'use', 'different']
```

```
100%|██████████| 4280/4280 [00:05<00:00, 846.96it/s]
```

```
(4280, 50)
[ 0.10992891 -0.13077738 -0.79036539  0.13625567  0.20219931 -0.0649115
1
 -0.19755909 -0.3391333  -0.0493638   0.13198713  0.18661448  0.0056173
4
 -0.05206201  0.47927658 -0.48073368 -0.36405726 -0.01985954  0.3774366
7
  0.057204    0.78528074 -0.64847556  0.57354722 -0.06889712 -0.7278922
4
  0.41433505 -0.15736638 -0.11427515 -0.18089402 -0.00195227  0.4773976
  0.1912275   0.7247466  -0.44553087  0.43363291 -0.37637368  0.1470013
4
  0.43232073 -0.1094127  -0.2849444  -0.15954012 -0.07255042 -0.3362244
8
  0.04545373  0.39030147  0.56404003 -0.24470929  0.49149002 -0.3319631
4
  0.07563705 -0.30794157]
```

**[5.2.1.4]tfidf-w2v**

In [0]:
```python
#this is for train data
i=0
list_of_sentance_train=[]
for sentance in X_train:
    list_of_sentance_train.append(sentance.split())


# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer()
tf_idf_matrix = model.fit_transform(X_train)
```

```python
# we are converting a dictionary with word as a key, and the idf as a v
alue
dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
```

```python
# TF-IDF weighted Word2Vec
tfidf_feat = model.get_feature_names() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and ce
ll_val = tfidf

tfidf_sent_vectors_train = []; # the tfidf-w2v for each sentence/review
 is stored in this list
row=0;
for sent in tqdm(list_of_sentance_train): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum =0; # num of words with a valid vector in the sentence/r
eview
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
#             tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    tfidf_sent_vectors_train.append(sent_vec)
    row += 1
tfidf_sent_vectors_train= np.array(sent_vectors_train)
print(tfidf_sent_vectors_train.shape)
print(tfidf_sent_vectors_train[0])
```

```
100%|██████████| 8687/8687 [01:28<00:00, 98.64it/s]
```

```
(8687, 50)
[ 0.21411499 -0.36776797 -0.75282662  0.33994114  0.40363468 -0.4364555
```

```
3
 -0.21620836 -0.39892652  0.31709211  0.3927634   0.31608421  0.2378958
1
  0.15208639  0.5907577  -0.59054357 -0.43573927  0.16540512  0.1558099
8
 -0.10617807  0.7075104  -0.6076003   0.54247047  0.00242196 -0.7233496
5
  0.53540209 -0.20284717 -0.0522027   0.04519267 -0.40957193  0.5115167
9
  0.29258774  0.76476256 -0.44577942  0.44159299 -0.44819606  0.3085005
3
  0.31808708  0.0550734  -0.35594978 -0.03504065 -0.07206956 -0.2751188
3
 -0.05680525  0.21498388  0.59119366  0.0805137   0.4885871  -0.5601444
2
 -0.00090877 -0.41330399]
```

In [0]:
```python
#this is for test data
i=0
list_of_sentance_test=[]
for sentance in X_test:
    list_of_sentance_test.append(sentance.split())


# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
#model = TfidfVectorizer()
tf_idf_matrix = model.transform(X_test)
# we are converting a dictionary with word as a key, and the idf as a v
alue
dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))



# TF-IDF weighted Word2Vec
tfidf_feat = model.get_feature_names() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and ce
ll_val = tfidf
```

```python
tfidf_sent_vectors_test = []; # the tfidf-w2v for each sentence/review
 is stored in this list
row=0;
for sent in tqdm(list_of_sentance_test): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum =0; # num of words with a valid vector in the sentence/r
eview
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
#              tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    tfidf_sent_vectors_test.append(sent_vec)
    row += 1
tfidf_sent_vectors_test= np.array(sent_vectors_test)
print(tfidf_sent_vectors_test.shape)
print(tfidf_sent_vectors_test[0])
```

```
100%|██████████| 6387/6387 [01:04<00:00, 98.43it/s]
```

```
(6387, 50)
[ 0.1020971  -0.25778366 -0.65248146  0.4684123   0.31574818 -0.4129766
5
 -0.20247924 -0.40664902  0.19689101  0.31480956  0.27975016  0.3556833
4
  0.26890199  0.6044926  -0.56634632 -0.39209652  0.14483991  0.0565501
4
  0.0810383   0.62992669 -0.66624729  0.55472585 -0.19655463 -0.7660524
2
  0.66199181 -0.15676928 -0.17650486 -0.0201654  -0.14755801  0.5113350
9
  0.20591084  0.8363934  -0.4064      0.41434614 -0.48395601  0.2026638
8
```

```
   0.45174919  0.05811068 -0.35995091  0.11331295 -0.02787627 -0.4291904
 1
 -0.06195572  0.35714837  0.65272359 -0.02911503  0.38616773 -0.5571473
  0.07343648 -0.32804483]
```

In [0]:
```python
#this is for cv data
i=0
list_of_sentance_cv=[]
for sentance in X_cv:
    list_of_sentance_cv.append(sentance.split())


# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
#model = TfidfVectorizer()
tf_idf_matrix = model.transform(X_cv)
# we are converting a dictionary with word as a key, and the idf as a v
alue
dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))



# TF-IDF weighted Word2Vec
tfidf_feat = model.get_feature_names() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and ce
ll_val = tfidf

tfidf_sent_vectors_cv = []; # the tfidf-w2v for each sentence/review is
 stored in this list
row=0;
for sent in tqdm(list_of_sentance_cv): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum =0; # num of words with a valid vector in the sentence/r
eview
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
#             tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
```

```python
            # dictionary[word] = idf value of word in whole courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    tfidf_sent_vectors_cv.append(sent_vec)
    row += 1
tfidf_sent_vectors_cv= np.array(sent_vectors_cv)
print(tfidf_sent_vectors_cv.shape)
print(tfidf_sent_vectors_cv[0])
```

```
100%|██████████| 4280/4280 [00:47<00:00, 89.49it/s]
```

```
(4280, 50)
[ 0.10992891 -0.13077738 -0.79036539  0.13625567  0.20219931 -0.0649115
 1
 -0.19755909 -0.3391333  -0.0493638   0.13198713  0.18661448  0.0056173
 4
 -0.05206201  0.47927658 -0.48073368 -0.36405726 -0.01985954  0.3774366
 7
  0.057204    0.78528074 -0.64847556  0.57354722 -0.06889712 -0.7278922
 4
  0.41433505 -0.15736638 -0.11427515 -0.18089402 -0.00195227  0.4773976
  0.1912275   0.7247466  -0.44553087  0.43363291 -0.37637368  0.1470013
 4
  0.43232073 -0.1094127  -0.2849444  -0.15954012 -0.07255042 -0.3362244
 8
  0.04545373  0.39030147  0.56404003 -0.24470929  0.49149002 -0.3319631
 4
  0.07563705 -0.30794157]
```

### [5.2.2] Applying RBF SVM on BOW, <span style="color:red">SET 1</span>

```python
In [0]: from sklearn.model_selection import train_test_split
        from sklearn.model_selection import learning_curve, GridSearchCV
        from sklearn.svm import SVC
```

```python
#Using GridSearchCV
tuned_parameters = [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
model = GridSearchCV(SVC(kernel='rbf'), tuned_parameters,scoring = 'roc
_auc', cv=5)

model.fit(X_train_bow, y_train)

print(model.best_estimator_)
print(model.score(X_test_bow, y_test))
```

```
SVC(C=100, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
    kernel='rbf', max_iter=-1, probability=False, random_state=None,
    shrinking=True, tol=0.001, verbose=False)
0.8869867599908448
```
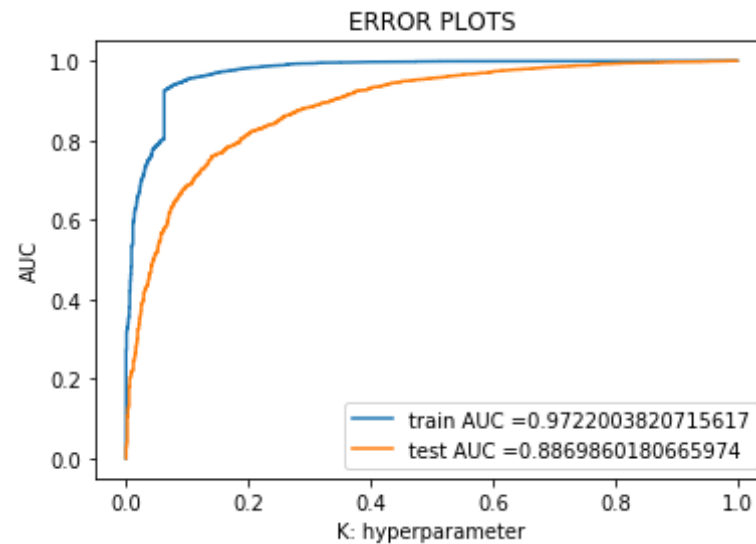
In [0]:
```python
from sklearn.svm import SVC
clf=SVC(C=100,kernel='rbf',probability=True)
clf.fit(X_train_bow,y_train)
train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba
(X_train_bow)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(X_
test_bow)[:,1])

plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, t
rain_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_
tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```

## ERROR PLOTS



```python
#for seaborn confusion matrix :https://stackoverflow.com/questions/3557
2000/how-can-i-plot-a-confusion-matrix
#refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
fusion-matrix-with-labels
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
ax= plt.subplot()
arr1=confusion_matrix(y_train, clf.predict(X_train_bow))
df_1= pd.DataFrame(arr1, range(2),range(2))
plt.figure(figsize = (5,2))
sn.heatmap(df_1, annot=True,fmt="d",ax=ax)
ax.set_title('Confusion Matrix');
ax.set_xlabel('Actual Labels')
ax.set_ylabel('Predicted Labels')
ax.xaxis.set_ticklabels(['False', 'True']);
ax.yaxis.set_ticklabels(['True', 'False']);
```

Train confusion matrix

Confusion Matrix

```
<Figure size 360x144 with 0 Axes>
```

In [0]: 
```python
#refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
fusion-matrix-with-labels
print("Test confusion matrix")
ax= plt.subplot()
arr2=confusion_matrix(y_test, clf.predict(X_test_bow))
df_2= pd.DataFrame(arr2, range(2),range(2))
plt.figure(figsize = (5,2))
sn.heatmap(df_2, annot=True,fmt="d",ax=ax)
ax.set_title('Confusion Matrix');
ax.set_xlabel('Actual Labels')
ax.set_ylabel('Predicted Labels')
ax.xaxis.set_ticklabels(['False', 'True']);
ax.yaxis.set_ticklabels(['True', 'False']);
```

Test confusion matrix

Confusion Matrix

```
<Figure size 360x144 with 0 Axes>
```

### [5.2.3] Applying RBF SVM on TFIDF, SET 2

In [0]:
```python
from sklearn.model_selection import train_test_split
from sklearn.model_selection import learning_curve, GridSearchCV
from sklearn.svm import SVC

#Using GridSearchCV
tuned_parameters = [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
model = GridSearchCV(SVC(kernel='rbf'), tuned_parameters,scoring = 'roc
_auc', cv=5)

model.fit(X_train_tf_idf, y_train)

print(model.best_estimator_)
print(model.score(X_test_tf_idf, y_test))
```
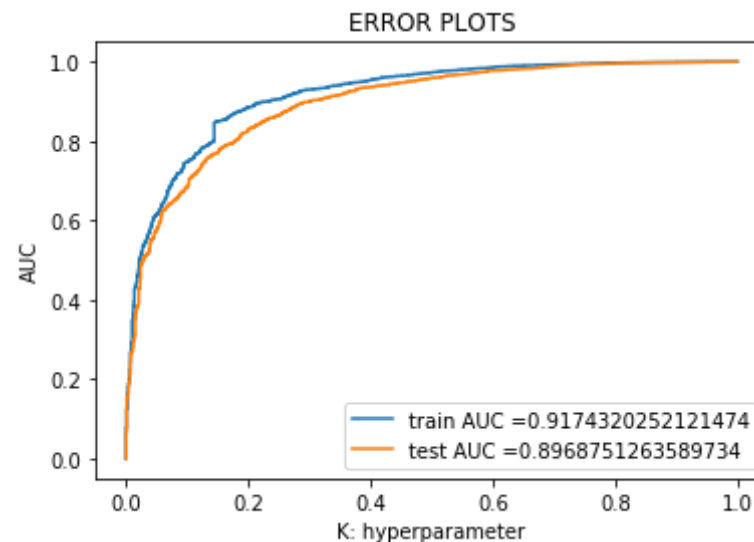
```
SVC(C=100, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
    kernel='rbf', max_iter=-1, probability=False, random_state=None,
```

```
              shrinking=True, tol=0.001, verbose=False)
          0.8968732715483552
```

In [0]:
```python
from sklearn.svm import SVC
clf=SVC(C=100,kernel='rbf',probability=True)
clf.fit(X_train_tf_idf,y_train)
train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba
(X_train_tf_idf)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(X_
test_tf_idf)[:,1])

plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, t
rain_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_
tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```



In [0]:
```python
#for seaborn confusion matrix :https://stackoverflow.com/questions/3557
```

```
2000/how-can-i-plot-a-confusion-matrix
#refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
fusion-matrix-with-labels
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
ax= plt.subplot()
arr1=confusion_matrix(y_train, clf.predict(X_train_tf_idf))
df_1= pd.DataFrame(arr1, range(2),range(2))
plt.figure(figsize = (5,2))
sn.heatmap(df_1, annot=True,fmt="d",ax=ax)
ax.set_title('Confusion Matrix');
ax.set_xlabel('Actual Labels')
ax.set_ylabel('Predicted Labels')
ax.xaxis.set_ticklabels(['False', 'True']);
ax.yaxis.set_ticklabels(['True', 'False']);
```

Train confusion matrix



Confusion Matrix

```
<Figure size 360x144 with 0 Axes>
```

```
In [0]:  #refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
         fusion-matrix-with-labels
         print("Test confusion matrix")
         ax= plt.subplot()
         arr2=confusion_matrix(y_test, clf.predict(X_test_tf_idf))
         df_2= pd.DataFrame(arr2, range(2),range(2))
         plt.figure(figsize = (5,2))
         sn.heatmap(df_2, annot=True,fmt="d",ax=ax)
         ax.set_title('Confusion Matrix');
         ax.set_xlabel('Actual Labels')
         ax.set_ylabel('Predicted Labels')
         ax.xaxis.set_ticklabels(['False', 'True']);
         ax.yaxis.set_ticklabels(['True', 'False']);
```

Test confusion matrix



<Figure size 360x144 with 0 Axes>

### [5.2.4] Applying RBF SVM on AVG W2V, SET 3

```
In [0]:  from sklearn.model_selection import train_test_split
```

```python
from sklearn.model_selection import learning_curve, GridSearchCV
from sklearn.svm import SVC

#Using GridSearchCV
tuned_parameters = [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
model = GridSearchCV(SVC(kernel='rbf'), tuned_parameters,scoring = 'roc
_auc', cv=5)

model.fit(sent_vectors_train, y_train)

print(model.best_estimator_)
print(model.score(sent_vectors_test, y_test))
```

```
SVC(C=100, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
    kernel='rbf', max_iter=-1, probability=False, random_state=None,
    shrinking=True, tol=0.001, verbose=False)
0.8540455088914056
```
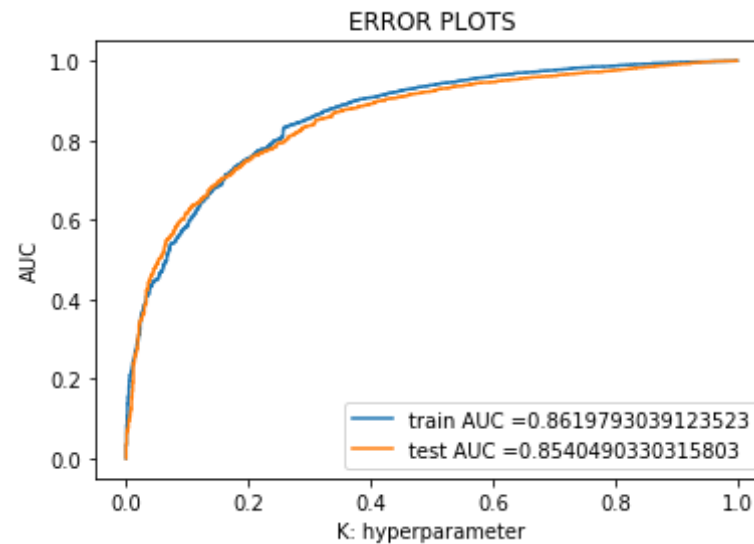
In [0]:
```python
from sklearn.svm import SVC
clf=SVC(C=100,kernel='rbf',probability=True)
clf.fit(sent_vectors_train,y_train)
train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba
(sent_vectors_train)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(se
nt_vectors_test)[:,1])

plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, t
rain_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_
tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```
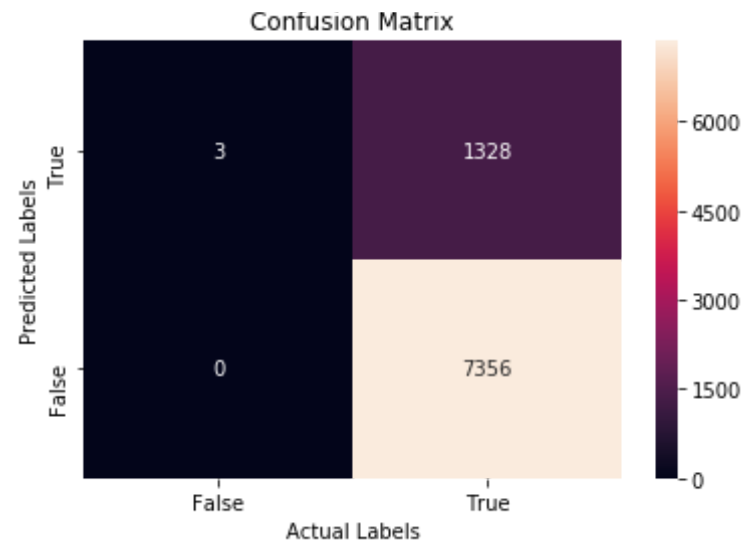
## ERROR PLOTS



```
In [0]:  #for seaborn confusion matrix :https://stackoverflow.com/questions/3557
         2000/how-can-i-plot-a-confusion-matrix
         #refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
         fusion-matrix-with-labels
         import seaborn as sn
         import pandas as pd
         import matplotlib.pyplot as plt
         from sklearn.metrics import confusion_matrix
         print("Train confusion matrix")
         ax= plt.subplot()
         arr1=confusion_matrix(y_train, clf.predict(sent_vectors_train))
         df_1= pd.DataFrame(arr1, range(2),range(2))
         plt.figure(figsize = (5,2))
         sn.heatmap(df_1, annot=True,fmt="d",ax=ax)
         ax.set_title('Confusion Matrix');
         ax.set_xlabel('Actual Labels')
         ax.set_ylabel('Predicted Labels')
         ax.xaxis.set_ticklabels(['False', 'True']);
         ax.yaxis.set_ticklabels(['True', 'False']);
```
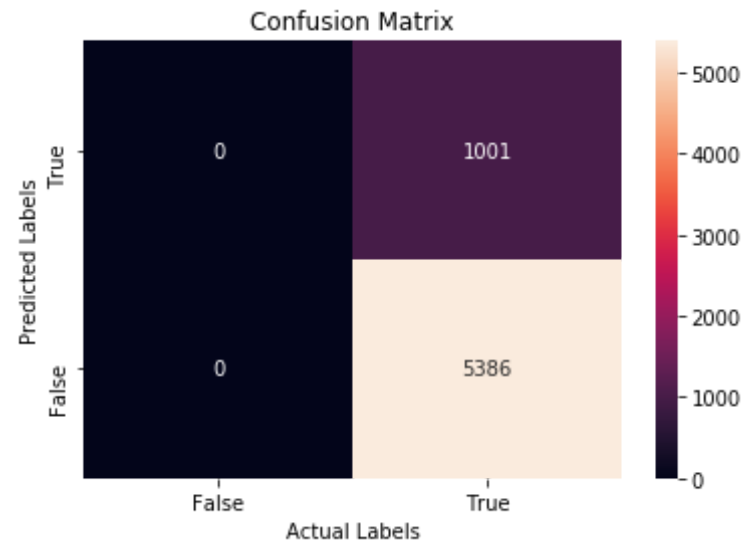
Train confusion matrix

Confusion Matrix

```
<Figure size 360x144 with 0 Axes>
```

In [0]:
```python
#refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
print("Test confusion matrix")
ax= plt.subplot()
arr2=confusion_matrix(y_test, clf.predict(sent_vectors_test))
df_2= pd.DataFrame(arr2, range(2),range(2))
plt.figure(figsize = (5,2))
sn.heatmap(df_2, annot=True,fmt="d",ax=ax)
ax.set_title('Confusion Matrix');
ax.set_xlabel('Actual Labels')
ax.set_ylabel('Predicted Labels')
ax.xaxis.set_ticklabels(['False', 'True']);
ax.yaxis.set_ticklabels(['True', 'False']);
```

Test confusion matrix

Confusion Matrix

```
<Figure size 360x144 with 0 Axes>
```

### [5.2.5] Applying RBF SVM on TFIDF W2V, <span style="color:red">SET 4</span>

In [0]:
```python
from sklearn.model_selection import train_test_split
from sklearn.model_selection import learning_curve, GridSearchCV
from sklearn.svm import SVC

#Using GridSearchCV
tuned_parameters = [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
model = GridSearchCV(SVC(kernel='rbf'), tuned_parameters,scoring = 'roc
_auc', cv=5)

model.fit(tfidf_sent_vectors_train, y_train)

print(model.best_estimator_)
print(model.score(tfidf_sent_vectors_test, y_test))
```
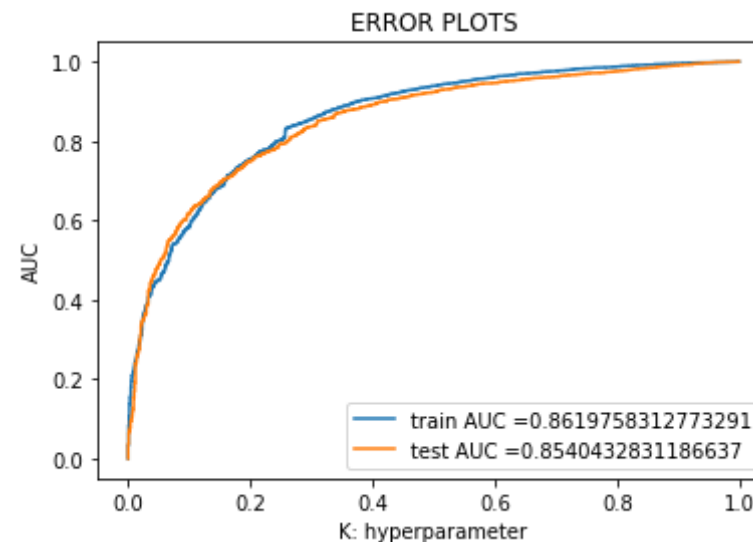
```
SVC(C=100, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
    kernel='rbf', max_iter=-1, probability=False, random_state=None,
```

```
        shrinking=True, tol=0.001, verbose=False)
        0.8540455088914056
```

In [0]:
```python
from sklearn.svm import SVC
clf=SVC(C=100,kernel='rbf',probability=True)
clf.fit(tfidf_sent_vectors_train,y_train)
train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba
(tfidf_sent_vectors_train)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(tf
idf_sent_vectors_test)[:,1])

plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, t
rain_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_
tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```
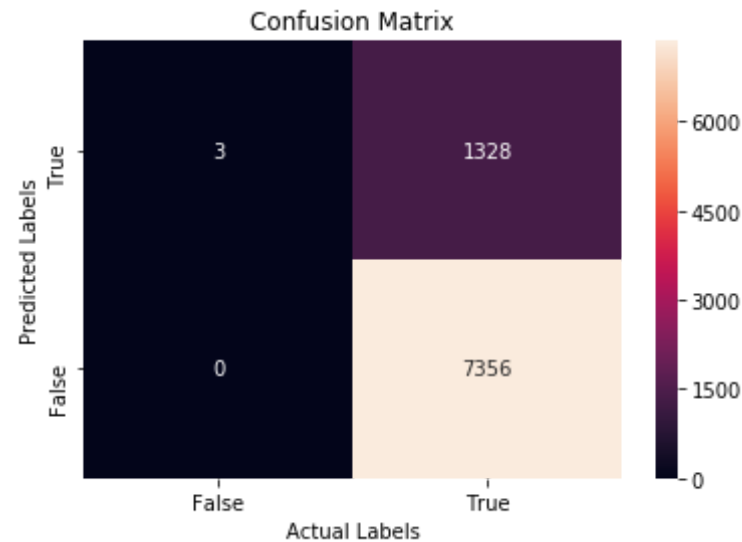


In [0]: `#for seaborn confusion matrix :https://stackoverflow.com/questions/3557`

```
2000/how-can-i-plot-a-confusion-matrix
#refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
fusion-matrix-with-labels
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
ax= plt.subplot()
arr1=confusion_matrix(y_train, clf.predict(tfidf_sent_vectors_train))
df_1= pd.DataFrame(arr1, range(2),range(2))
plt.figure(figsize = (5,2))
sn.heatmap(df_1, annot=True,fmt="d",ax=ax)
ax.set_title('Confusion Matrix');
ax.set_xlabel('Actual Labels')
ax.set_ylabel('Predicted Labels')
ax.xaxis.set_ticklabels(['False', 'True']);
ax.yaxis.set_ticklabels(['True', 'False']);
```

Train confusion matrix



<Figure size 360x144 with 0 Axes>

```
In [0]: #refernce:https://stackoverflow.com/questions/19233771/sklearn-plot-con
        fusion-matrix-with-labels
        print("Test confusion matrix")
        ax= plt.subplot()
        arr2=confusion_matrix(y_test, clf.predict(sent_vectors_test))
        df_2= pd.DataFrame(arr2, range(2),range(2))
        plt.figure(figsize = (5,2))
        sn.heatmap(df_2, annot=True,fmt="d",ax=ax)
        ax.set_title('Confusion Matrix');
        ax.set_xlabel('Actual Labels')
        ax.set_ylabel('Predicted Labels')
        ax.xaxis.set_ticklabels(['False', 'True']);
        ax.yaxis.set_ticklabels(['True', 'False']);
```
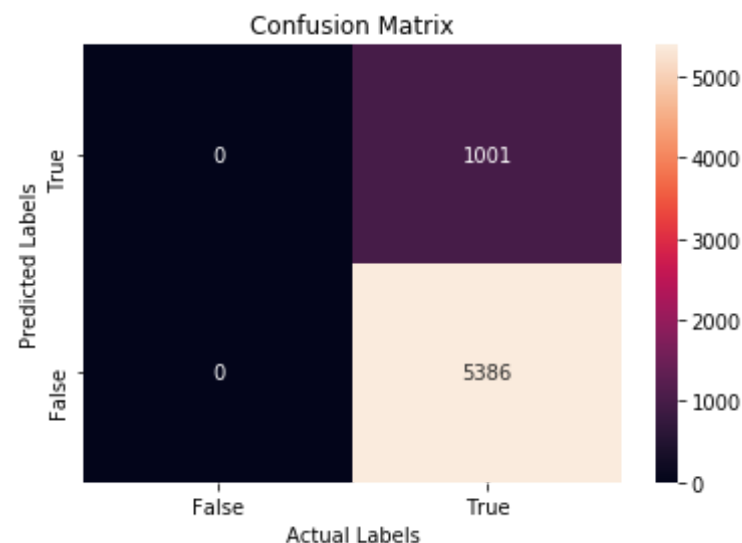
Test confusion matrix



`<Figure size 360x144 with 0 Axes>`

## [6] Conclusions

```
In [118]: from prettytable import PrettyTable
```

```
print('auc performace table:')
x = PrettyTable()
x.field_names = ["Vectorizer", "Regularizer", 'alpha',"auc"]
x.add_row(["BoW", "L2", 0.01,0.92])
x.add_row(["tfidf", "L2",0.0001, 0.95])
x.add_row(["avg w2v", "L2", 0.01,0.88])
x.add_row(["tfidfw2v", "L2",0.01, 0.89])
print(x)
y = PrettyTable()
y.field_names = ["Vectorizer",'C',"auc"]
y.add_row(["BoW",100,0.88])
y.add_row(["tfidf",100,0.89])
y.add_row(["avg w2v",100,0.85])
y.add_row(["tfidf",100,0.85])
print(y)
```

```
auc performace table:
+------------+-------------+--------+------+
| Vectorizer | Regularizer | alpha  | auc  |
+------------+-------------+--------+------+
|    BoW     |      L2     |  0.01  | 0.92 |
|   tfidf    |      L2     | 0.0001 | 0.95 |
|  avg w2v   |      L2     |  0.01  | 0.88 |
|  tfidfw2v  |      L2     |  0.01  | 0.89 |
+------------+-------------+--------+------+
+------------+-----+------+
| Vectorizer |  C  | auc  |
+------------+-----+------+
|    BoW     | 100 | 0.88 |
|   tfidf    | 100 | 0.89 |
|  avg w2v   | 100 | 0.85 |
|   tfidf    | 100 | 0.85 |
+------------+-----+------+
```