# 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

## [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 [1]:  %matplotlib inline
         import warnings
         warnings.filterwarnings("ignore")


         import sqlite3
         import pandas as pd
         import numpy as np
         import nltk
         import string
```

```python
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
```

```
C:\Users\LiGht\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWa
rning: detected Windows; aliasing chunkize to chunkize_serial
  warnings.warn("detected Windows; aliasing chunkize to chunkize_seria
l")
```

In [2]:
```python
# using SQLite Table to read data.
con = sqlite3.connect('database.sqlite')

# filtering only positive and negative reviews i.e.
# not taking into consideration those reviews with Score=3
# SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000, will give top 50
0000 data points
# you can change the number to any other number based on your computing
 power
```

```python
# filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Sco
re != 3 LIMIT 60000""", con)
# for tsne assignment you can take 5k data points

filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score
 != 3 LIMIT 60000""", con)

# Give reviews with Score>3 a positive rating(1), and reviews with a sc
ore<3 a negative rating(0).
def partition(x):
    if x < 3:
        return 0
    return 1

#changing reviews with score less than 3 to be positive and vice-versa
actualScore = filtered_data['Score']
positiveNegative = actualScore.map(partition)
filtered_data['Score'] = positiveNegative
print("Number of data points in our data", filtered_data.shape)
filtered_data.head(3)
```

Number of data points in our data (60000, 10)

Out[2]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfulnes |
|---|---|---|---|---|---|---|
| 0 | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1 | 1 |
| 1 | 2 | B00813GRG4 | A1D87F6ZCVE5NK | dll pa | 0 | 0 |

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfulnes |
|---|---|---|---|---|---|---|
| **2** | 3 | B000LQOCH0 | ABXLMWJIXXAIN | Natalia Corres "Natalia Corres" | 1 | 1 |

```
In [3]: display = pd.read_sql_query("""
        SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
        FROM Reviews
        GROUP BY UserId
        HAVING COUNT(*)>1
        """, con)
```

```
In [4]: print(display.shape)
        display.head()
```

(80668, 7)

Out[4]:

| | UserId | ProductId | ProfileName | Time | Score | Text | COUI |
|---|---|---|---|---|---|---|---|
| **0** | #oc-R115TNMSPFT9I7 | B007Y59HVM | Breyton | 1331510400 | 2 | Overall its just OK when considering the price... | 2 |

| | UserId | ProductId | ProfileName | Time | Score | Text | COU |
|---|---|---|---|---|---|---|---|
| 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 [5]: `display[display['UserId']=='AZY10LLTJ71NX']`

Out[5]:

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

```
In [6]: display['COUNT(*)'].sum()
```

Out[6]: 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 [7]: display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID
""", con)
display.head()
```

Out[7]:

|   | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfuln |
|---|------|-----------|---------------|------------------|----------------------|----------|
| 0 | 78445 | B000HDL1RQ | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |

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

```
In [8]: #Sorting data according to ProductId in ascending order
        sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=Tr
        ue, inplace=False, kind='quicksort', na_position='last')
```

```
In [9]: #Deduplication of entries
        final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time"
        ,"Text"}, keep='first', inplace=False)
        final.shape
```

```
Out[9]: (54458, 10)
```

```
In [10]: #Checking to see how much % of data still remains
         (final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

```
Out[10]: 90.76333333333334
```

**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 [11]: display= pd.read_sql_query("""
         SELECT *
         FROM Reviews
```

```
WHERE Score != 3 AND Id=44737 OR Id=64422
ORDER BY ProductID
""", con)

display.head()
```

Out[11]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfuln |
|---|---|---|---|---|---|---|
| 0 | 64422 | B000MIDROQ | A161DK06JJMCYF | J. E. Stephens "Jeanne" | 3 | 1 |
| 1 | 44737 | B001EQ55RW | A2V0I904FH7ABY | Ram | 3 | 2 |

In [12]:
```
final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

In [13]:
```
#Before starting the next phase of preprocessing lets see the number of
 entries left
print(final.shape)

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

```
(54456, 10)
```

Out[13]:
```
1    45572
0     8884
```

```
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 [14]:  # printing some random reviews
          sent_0 = final['Text'].values[0]
          print(sent_0)
          print("="*50)

          sent_1000 = final['Text'].values[1000]
          print(sent_1000)
          print("="*50)

          sent_1500 = final['Text'].values[1500]
          print(sent_1500)
```

```python
print("="*50)

sent_4900 = final['Text'].values[4900]
print(sent_4900)
print("="*50)
```

```
Our dogs just love them.  I saw them in a pet store and a tag was attac
hed regarding them being made in China and it satisfied me that they we
re safe.
==================================================
My whole family loves this sauce. I use it primarily to recreate a dish
we had in Maui. It adds a wonderful sweet and spicy flavor to just abou
t any sauce. Service on this order was spectacular. Several bottles arr
ived broken and replacements were received within days. Amazon service
was fast, easy and reliable. I love Amazon. Also, there are many recipe
s that you can look up using this sauce that a very good. It goes spect
aculary with coconut milk, vegetables, chicken or shrimp and pasta. We
love it.
==================================================
My husband is a paraplegic and was having UTIs constantly.  Since he st
arted drinking a cup of this tea every morning, he hasn't had a UTI for
a year.
==================================================
THIS BREAD MIX IS THE CLOSEST THING TO REGULAR BREAD I HAVE FOUND. EASY
TO MAKE IN MY BREAD MACHINE.
==================================================
```

In [15]:
```python
# remove urls from text python: https://stackoverflow.com/a/40823105/40
84039
sent_0 = re.sub(r"http\S+", "", sent_0)
sent_1000 = re.sub(r"http\S+", "", sent_1000)
sent_150 = re.sub(r"http\S+", "", sent_1500)
sent_4900 = re.sub(r"http\S+", "", sent_4900)

print(sent_0)
```

```
Our dogs just love them.  I saw them in a pet store and a tag was attac
hed regarding them being made in China and it satisfied me that they we
re safe.
```

```python
In [16]:  # https://stackoverflow.com/questions/16206380/python-beautifulsoup-how
          -to-remove-all-tags-from-an-element
          from bs4 import BeautifulSoup

          soup = BeautifulSoup(sent_0, 'lxml')
          text = soup.get_text()
          print(text)
          print("="*50)

          soup = BeautifulSoup(sent_1000, 'lxml')
          text = soup.get_text()
          print(text)
          print("="*50)

          soup = BeautifulSoup(sent_1500, 'lxml')
          text = soup.get_text()
          print(text)
          print("="*50)

          soup = BeautifulSoup(sent_4900, 'lxml')
          text = soup.get_text()
          print(text)
```

```
Our dogs just love them.  I saw them in a pet store and a tag was attac
hed regarding them being made in China and it satisfied me that they we
re safe.
==================================================
My whole family loves this sauce. I use it primarily to recreate a dish
we had in Maui. It adds a wonderful sweet and spicy flavor to just abou
t any sauce. Service on this order was spectacular. Several bottles arr
ived broken and replacements were received within days. Amazon service
was fast, easy and reliable. I love Amazon. Also, there are many recipe
s that you can look up using this sauce that a very good. It goes spect
aculary with coconut milk, vegetables, chicken or shrimp and pasta. We
love it.
==================================================
My husband is a paraplegic and was having UTIs constantly.  Since he st
arted drinking a cup of this tea every morning, he hasn't had a UTI for
a year.
```

```
==================================================
THIS BREAD MIX IS THE CLOSEST THING TO REGULAR BREAD I HAVE FOUND. EASY
TO MAKE IN MY BREAD MACHINE.
```

In [17]:
```python
# 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
    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 [18]:
```python
sent_1500 = decontracted(sent_1500)
print(sent_1500)
print("="*50)
```

```
My husband is a paraplegic and was having UTIs constantly.  Since he st
arted drinking a cup of this tea every morning, he has not had a UTI fo
r a year.
==================================================
```

In [19]:
```python
#remove words with numbers python: https://stackoverflow.com/a/1808237
0/4084039
sent_0 = re.sub("\S*\d\S*", "", sent_0).strip()
print(sent_0)
```

```
Our dogs just love them.  I saw them in a pet store and a tag was attac
hed regarding them being made in China and it satisfied me that they we
```

re safe.

In [20]:
```python
#remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)
print(sent_1500)
```

My husband is a paraplegic and was having UTIs constantly Since he star
ted drinking a cup of this tea every morning he has not had a UTI for a
year

In [21]:
```python
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'no
t'
# <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', 'o
urs', 'ourselves', 'you', "you're", "you've",\
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselve
s', 'he', 'him', 'his', 'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'it
s', 'itself', 'they', 'them', 'their',\
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'th
is', 'that', "that'll", 'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'h
ave', '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', 'h
ow', 'all', 'any', 'both', 'each', 'few', 'more',\
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 's
o', 'than', 'too', 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should',
```

```
             "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 [22]:
```python
# Combining all the above stundents
from tqdm import tqdm
preprocessed_reviews = []
# tqdm is for printing the status bar
for sentance in tqdm(final['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.append(sentance.strip())
```

```
100%|████████████████████████████████████████████████████████████████
████████| 54456/54456 [00:26<00:00, 2023.24it/s]
```

In [23]:
```python
preprocessed_reviews[1500]
```

Out[23]:
```
'husband paraplegic utis constantly since started drinking cup tea ever
y morning not uti year'
```

## [3.2] Preprocessing Review Summary

In [24]:
```python
## Similartly you can do preprocessing for review summary also.
```

# [4] Featurization

## [4.1] BAG OF WORDS

```
In [25]: #BoW
         count_vect = CountVectorizer() #in scikit-learn
         count_vect.fit(preprocessed_reviews)
         print("some feature names ", count_vect.get_feature_names()[:10])
         print('='*50)

         final_counts = count_vect.transform(preprocessed_reviews)
         print("the type of count vectorizer ",type(final_counts))
         print("the shape of out text BOW vectorizer ",final_counts.get_shape())
         print("the number of unique words ", final_counts.get_shape()[1])
```

```
some feature names  ['aa', 'aaa', 'aaaa', 'aaaaa', 'aaaaaaaaaaaa', 'aaa
aaaaaaaaaaa', 'aaaaaaahhhhhh', 'aaaaaawwwwwwwww', 'aaaaah', 'aaaand']
==================================================
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (54456, 43297)
the number of unique words  43297
```

## [4.2] Bi-Grams and n-Grams.

```
In [26]: #bi-gram, tri-gram and n-gram

         #removing stop words like "not" should be avoided before building n-gra
         ms
         # count_vect = CountVectorizer(ngram_range=(1,2))
         # please do read the CountVectorizer documentation http://scikit-learn.
         org/stable/modules/generated/sklearn.feature_extraction.text.CountVecto
         rizer.html

         # you can choose these numebrs min_df=10, max_features=5000, of your ch
         oice
```

```python
count_vect = CountVectorizer(ngram_range=(1,2), min_df=10, max_features
=5000)
final_bigram_counts = count_vect.fit_transform(preprocessed_reviews)
print("the type of count vectorizer ",type(final_bigram_counts))
print("the shape of out text BOW vectorizer ",final_bigram_counts.get_s
hape())
print("the number of unique words including both unigrams and bigrams "
, final_bigram_counts.get_shape()[1])
```

```
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (54456, 5000)
the number of unique words including both unigrams and bigrams  5000
```

## [4.3] TF-IDF

In [27]:
```python
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
tf_idf_vect.fit(preprocessed_reviews)
print("some sample features(unique words in the corpus)",tf_idf_vect.ge
t_feature_names()[0:10])
print('='*50)

final_tf_idf = tf_idf_vect.transform(preprocessed_reviews)
print("the type of count vectorizer ",type(final_tf_idf))
print("the shape of out text TFIDF vectorizer ",final_tf_idf.get_shape
())
print("the number of unique words including both unigrams and bigrams "
, final_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 fe
ed', 'able figure', 'able find']
==================================================
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer  (54456, 32341)
the number of unique words including both unigrams and bigrams  32341
```

## [4.4] Word2Vec

In [28]:
```python
# Train your own Word2Vec model using your own text corpus
i=0
list_of_sentance=[]
for sentance in preprocessed_reviews:
    list_of_sentance.append(sentance.split())
```

In [29]:
```python
# Using Google News Word2Vectors

# in this project we are using a pretrained model by google
# its 3.3G file, once you load this into your memory
# it occupies ~9Gb, so please do this step only if you have >12G of ram
# we will provide a pickle file wich contains a dict ,
# and it contains all our courpus words as keys and  model[word] as values
# To use this code-snippet, download "GoogleNews-vectors-negative300.bin"
# from https://drive.google.com/file/d/0B7XkCwpI5KDYNlNUTTlSS21pQmM/edit
# it's 1.9GB in size.


# http://kavita-ganesan.com/gensim-word2vec-tutorial-starter-code/#.W17SRFAzZPY
# you can comment this whole cell
# or change these varible according to your need

is_your_ram_gt_16g=False
want_to_use_google_w2v = False
want_to_train_w2v = True

if want_to_train_w2v:
    # min_count = 5 considers only words that occured atleast 5 times
    w2v_model=Word2Vec(list_of_sentance,min_count=5,size=50, workers=4)
    print(w2v_model.wv.most_similar('great'))
    print('='*50)
    print(w2v_model.wv.most_similar('worst'))
```

```python
elif want_to_use_google_w2v and is_your_ram_gt_16g:
    if os.path.isfile('GoogleNews-vectors-negative300.bin'):
        w2v_model=KeyedVectors.load_word2vec_format('GoogleNews-vectors
-negative300.bin', binary=True)
        print(w2v_model.wv.most_similar('great'))
        print(w2v_model.wv.most_similar('worst'))
    else:
        print("you don't have gogole's word2vec file, keep want_to_trai
n_w2v = True, to train your own w2v ")
```

```
[('fantastic', 0.8443485498428345), ('awesome', 0.8417231440544128),
('good', 0.8162604570388794), ('amazing', 0.7886624932289124), ('excell
ent', 0.7871252298355103), ('terrific', 0.7861415147781372), ('perfec
t', 0.7453365325927734), ('wonderful', 0.7384995222091675), ('decent',
0.7292306423187256), ('fabulous', 0.6922100782394409)]
==================================================
[('greatest', 0.7277129292488098), ('best', 0.7236092686653137), ('tast
iest', 0.7223793268203735), ('closest', 0.6567206978797913), ('experien
ced', 0.6420025825500488), ('horrible', 0.6341232061386108), ('nasties
t', 0.6291502714157104), ('disgusting', 0.6223399639129639), ('ive', 0.
6210818290710449), ('awful', 0.615583062171936)]
```

In [30]:
```python
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])
```
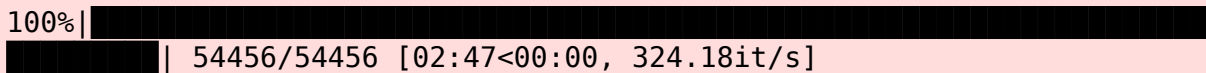
```
number of words that occured minimum 5 times  13942
sample words  ['dogs', 'love', 'saw', 'pet', 'store', 'tag', 'attache
d', 'regarding', 'made', 'china', 'satisfied', 'safe', 'loves', 'chicke
n', 'product', 'wont', 'buying', 'anymore', 'hard', 'find', 'products',
'usa', 'one', 'isnt', 'bad', 'good', 'take', 'chances', 'till', 'know',
'going', 'imports', 'available', 'victor', 'traps', 'unreal', 'course',
'total', 'fly', 'pretty', 'stinky', 'right', 'nearby', 'used', 'bait',
'seasons', 'ca', 'not', 'beat', 'great']
```

## [4.4.1] Converting text into vectors using Avg W2V, TFIDF-W2V

**[4.4.1.1] Avg W2v**

In [31]:
```python
# average Word2Vec
# compute average word2vec for each review.
sent_vectors = []; # the avg-w2v for each sentence/review is stored in
 this list
for sent in tqdm(list_of_sentance): # 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.append(sent_vec)
print(len(sent_vectors))
print(len(sent_vectors[0]))
```

```
100%|██████████████████████████████████████████████████████████████| 54456/54456 [02:47<00:00, 324.18it/s]
54456
50
```

**[4.4.1.2] TFIDF weighted W2v**

In [32]:
```python
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer()
tf_idf_matrix = model.fit_transform(preprocessed_reviews)
# 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_)))
```

```
In [33]: # 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 = []; # the tfidf-w2v for each sentence/review is st
         ored in this list
         row=0;
         for sent in tqdm(list_of_sentance): # 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.append(sent_vec)
             row += 1
```

```
100%|████████████████████████████████████████████████████████████████
████████| 54456/54456 [30:44<00:00, 29.53it/s]
```

# [5] Assignment 3: KNN

1. **Apply Knn(brute force version) 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. **Apply Knn(kd tree version) on these feature sets**
   NOTE: sklearn implementation of kd-tree accepts only dense matrices, you need to convert the sparse matrices of CountVectorizer/TfidfVectorizer into dense matices. You can convert sparse matrices to dense using .toarray() attribute. For more information please visit this link

   - SET 5:Review text, preprocessed one converted into vectors using (BOW) but with restriction on maximum features generated.

     ```
     count_vect = CountVectorizer(min_df=10,
     max_features=500)
     count_vect.fit(preprocessed_reviews)
     ```

   - SET 6:Review text, preprocessed one converted into vectors using (TFIDF) but with restriction on maximum features generated.

     ```
     tf_idf_vect = TfidfVectorizer(min_d
     f=10, max_features=500)
     tf_idf_vect.fit(preprocessed_review
     s)
     ```

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

3. **The hyper paramter tuning(find best K)**

   - 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. **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](confusion matrix) with predicted and original labels of test data points

  

5. **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](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.](link.)

# [5.1] Applying KNN brute force

**[5.1.1] Applying KNN brute force on BOW, SET 1**

In [30]: 
```python
# Please write all the code with proper documentation
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")
```

In [31]: 
```python
from sklearn import datasets, neighbors
final['CleanedText']=preprocessed_reviews
print(final.shape)
final.head(3)
```

(54456, 11)

Out[31]:

|  | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Help |
|---|---|---|---|---|---|---|
| **22621** | 24751 | 2734888454 | A1C298ITT645B6 | Hugh G. Pritchard | 0 | 0 |
| **22620** | 24750 | 2734888454 | A13ISQV0U9GZIC | Sandikaye | 1 | 1 |

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Help |
|---|---|---|---|---|---|---|
| **2546** | 2774 | B00002NCJC | A196AJHU9EASJN | Alex Chaffee | 0 | 0 |

In [32]:
```python
labels=final['Score'].values

X = final['CleanedText'].values
Y = labels[0:60000]
```

In [33]:
```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


print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer()
vectorizer.fit(X_train) # fit has to happen only on train data

# 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)
```

```
(24444,) (24444,)
(12041,) (12041,)
(17971,) (17971,)
After vectorizations
(24444, 29619) (24444,)
(12041, 29619) (12041,)
(17971, 29619) (17971,)
========================================================================
============================
```
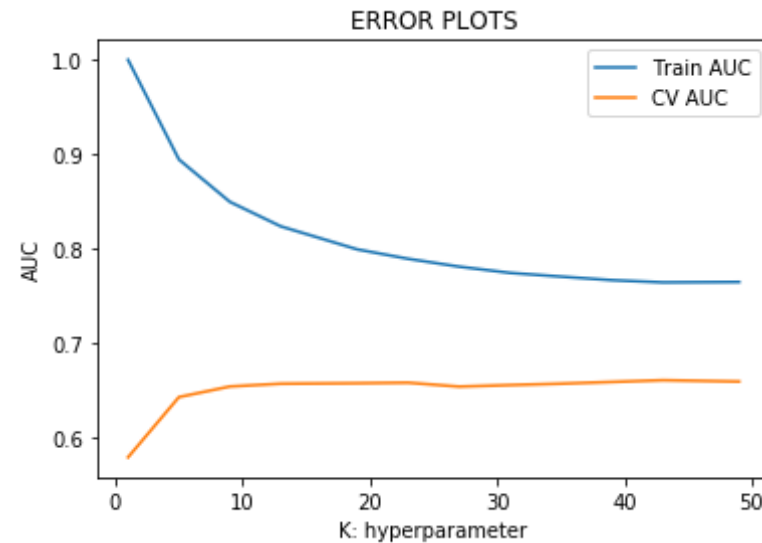
In [34]:
```
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score
from sklearn.cross_validation import cross_val_score
from collections import Counter
from sklearn import cross_validation

train_auc = []
cv_auc = []
K = [1,5,9,13,19,23,27,31,35,39,43,49]
for i in K:
    neigh = KNeighborsClassifier(n_neighbors=i,algorithm='brute')
    neigh.fit(X_train_bow, y_train)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probab
ility estimates of the positive class
    # not the predicted outputs
    y_train_pred =  neigh.predict_proba(X_train_bow)[:,1]
    y_cv_pred =  neigh.predict_proba(X_cv_bow)[:,1]

    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
```

```
plt.plot(K, train_auc, label='Train AUC')
plt.plot(K, cv_auc, label='CV AUC')
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```



In [35]:
```
#here we are choosing the best_k based on AUC result
best_k = 31
```

In [40]:
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc
_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc


neigh = KNeighborsClassifier(n_neighbors=best_k,algorithm='brute')
neigh.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
```
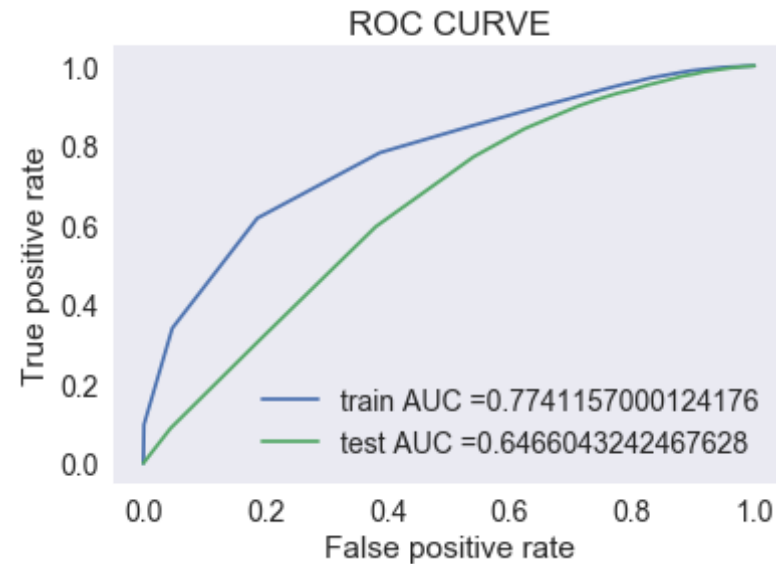
```
# not the predicted outputs

y_train_pred = neigh.predict_proba(X_train_bow)[:,1]
y_test_pred = neigh.predict_proba(X_test_bow)[:,1]

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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("False positive rate")
plt.ylabel("True positive rate")
plt.title("ROC CURVE")
plt.grid()
plt.show()
```



```
In [37]:  # we will pick a threshold that will give the least fpr
          def find_best_threshold(threshould, fpr, tpr):
```

```
        t = threshould[np.argmax(tpr*(1-fpr))]
        # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is
     very high
        print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for th
     reshold", np.round(t,3))
        return t

    def predict_with_best_t(proba, threshould):
        predictions = []
        for i in proba:
            if i>=threshould:
                predictions.append(1)
            else:
                predictions.append(0)
        return predictions
```

In [38]:
```
    print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_
t)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t
)))
```

```
========================================================================
============================
the maximum value of tpr*(1-fpr) 0.5016697263975959 for threshold 0.935
Train confusion matrix
[[ 3210   736]
 [ 7857 12641]]
Test confusion matrix
[[1828 1123]
 [6075 8945]]
```
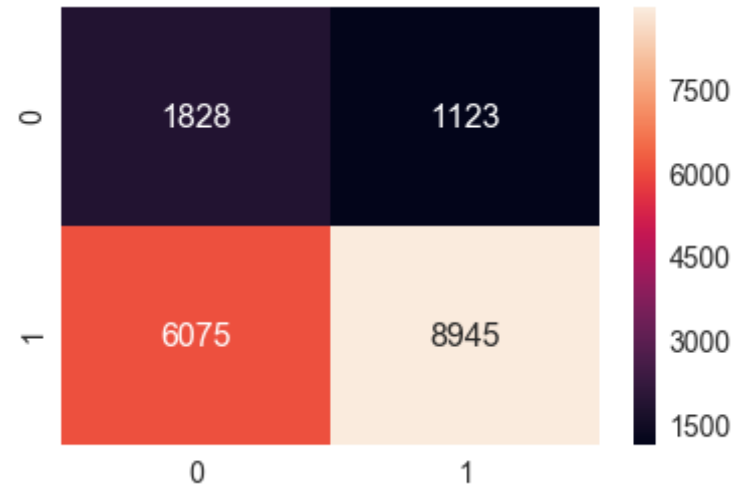
In [39]:
```
df_cm = pd.DataFrame((confusion_matrix(y_test, predict_with_best_t(y_te
st_pred, best_t))), range(2),range(2))
```

```
sns.set(font_scale=1.4)#for label size
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
```

Out[39]: `<matplotlib.axes._subplots.AxesSubplot at 0x5267a11518>`



**[5.1.2] Applying KNN brute force on TFIDF, <span style="color:red">SET 2</span>**

In [41]:
```
# Please write all the code with proper documentation
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


print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

from sklearn.feature_extraction.text import TfidfVectorizer
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
tf_idf_vect.fit(X_train) # fit has to happen only on train data
```

```python
# we use the fitted CountVectorizer to convert the text to vector
X_train_tfidf = tf_idf_vect.transform(X_train)
X_cv_tfidf = tf_idf_vect.transform(X_cv)
X_test_tfidf = tf_idf_vect.transform(X_test)

print("After vectorizations")
print(X_train_tfidf.shape, y_train.shape)
print(X_cv_tfidf.shape, y_cv.shape)
print(X_test_tfidf.shape, y_test.shape)
print("="*100)
```

```
(24444,) (24444,)
(12041,) (12041,)
(17971,) (17971,)
After vectorizations
(24444, 14620) (24444,)
(12041, 14620) (12041,)
(17971, 14620) (17971,)
================================================================================
===========================
```

```python
In [42]: train_auc = []
cv_auc = []
K = [1,5,9,13,19,23,27,31,35,39,43,49]
for i in K:
    neigh = KNeighborsClassifier(n_neighbors=i,algorithm='brute')
    neigh.fit(X_train_tfidf, y_train)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probab
ility estimates of the positive class
    # not the predicted outputs
    y_train_pred =  neigh.predict_proba(X_train_tfidf)[:,1]
    y_cv_pred =  neigh.predict_proba(X_cv_tfidf)[:,1]

    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(K, train_auc, label='Train AUC')
plt.plot(K, cv_auc, label='CV AUC')
plt.legend()
```
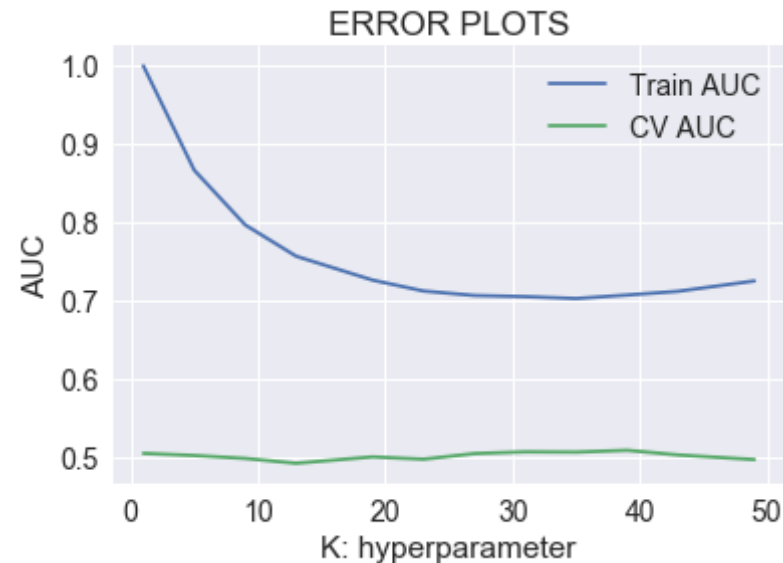
```python
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```



```
In [46]:  #here we are choosing the best_k based on AUC results
          best_k = 25
```

```python
In [47]:  from sklearn.metrics import roc_curve, auc


          neigh = KNeighborsClassifier(n_neighbors=best_k,algorithm='brute')
          neigh.fit(X_train_tfidf, 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

          y_train_pred = neigh.predict_proba(X_train_tfidf)[:,1]
          y_test_pred = neigh.predict_proba(X_test_tfidf)[:,1]

          train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
```
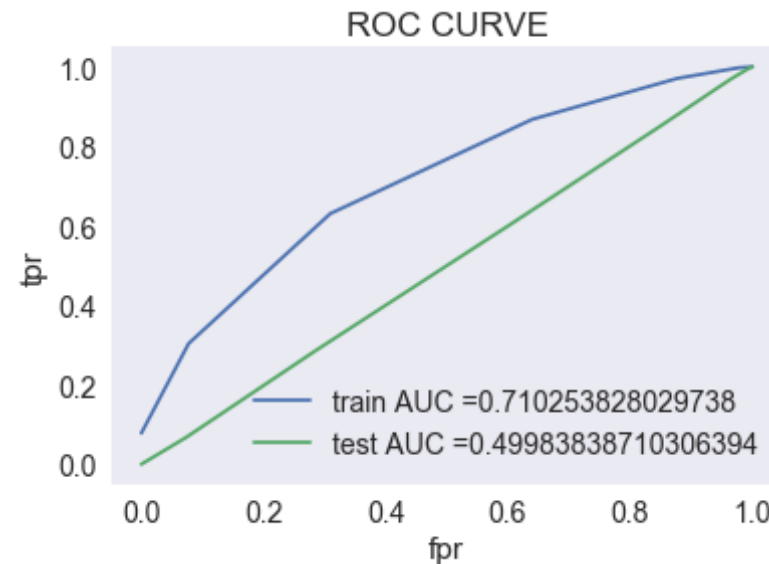
```python
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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("fpr")
plt.ylabel("tpr")
plt.title("ROC CURVE")
plt.grid()
plt.show()
```

ROC CURVE

train AUC =0.710253828029738
test AUC =0.49983838710306394

In [48]:
```python
print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_
t)))
print("Test confusion matrix")
```

```
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t
)))
```

```
=========================================================================
==============================
the maximum value of tpr*(1-fpr) 0.4356509144765014 for threshold 0.92
Train confusion matrix
[[ 2746  1230]
 [ 7557 12911]]
Test confusion matrix
[[1189 1778]
 [6020 8984]]
```
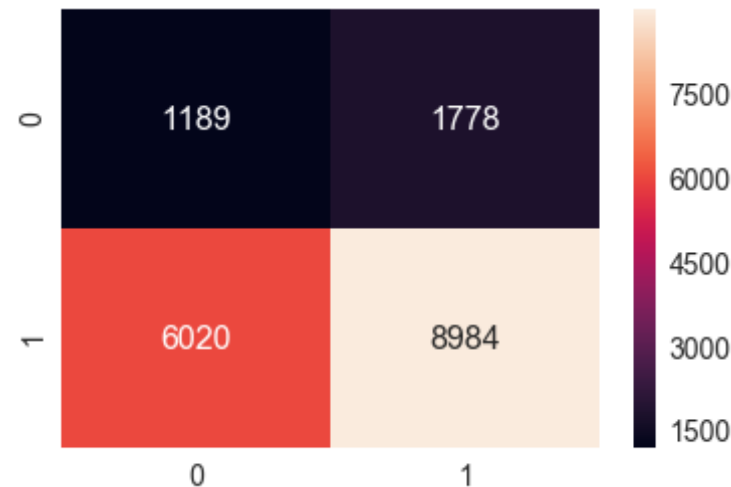
In [49]:
```
df_cm = pd.DataFrame((confusion_matrix(y_test, predict_with_best_t(y_te
st_pred, best_t))), range(2),range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
```

Out[49]: `<matplotlib.axes._subplots.AxesSubplot at 0x526aebf940>`



**[5.1.3] Applying KNN brute force on AVG W2V, SET 3**

```python
In [50]:  # Please write all the code with proper documentation
          i=0
          list_of_sentance_train=[]
          for sentance in X_train:
              list_of_sentance_train.append(sentance.split())
```

```python
In [51]:  # 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)
```

```python
In [52]:  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])
```

```
number of words that occured minimum 5 times  9523
sample words  ['thoroughly', 'enjoyed', 'buying', 'coffee', 'starbuck
s', 'lover', 'never', 'seen', 'line', 'curious', 'definitely', 'plus',
'received', 'items', 'initial', 'shipped', 'date', 'delivered', 'fairl
y', 'quickly', 'well', 'would', 'buy', 'favor', 'price', 'steal', 'oh',
'gluten', 'free', 'pantry', 'mess', 'great', 'thing', 'reviews', 'produ
ct', 'pre', 'recipe', 'change', 'wonder', 'everyone', 'thinks', 'change
d', 'first', 'ingredient', 'brownies', 'used', 'dutch', 'cocoa', 'decad
ent', 'not']
```

```python
In [53]:  # 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])
```

```
100%|████████████████████████████████████████████
████████| 24444/24444 [00:47<00:00, 519.93it/s]
```

```
(24444, 50)
[-0.04975627  0.06066158 -0.22982676 -0.00227746  0.06773219  0.3492490
9
 -0.07361499  0.69708604 -0.57257522  0.31711277  0.25784669  0.3887331
6
 -1.17572511  0.59886487 -0.51252999 -0.89382662 -0.45467491  0.7647610
8
 -0.28513595  0.06017749 -0.46121079  0.22203715 -0.48193237  0.0322674
1
  0.13213545 -0.59774461  0.43969732  0.41556463  0.05049974  0.3505374
4
 -0.23647515 -0.91181191  0.28498196  0.20754382  0.57985938  0.5128986
5
  0.76918463 -0.32573557  0.229941    0.02619685 -0.949994   -0.2571002
1
  0.53441468 -0.3976952  -0.64212655 -0.15986535  0.16180914 -0.1594124
 -0.00248328 -0.58835612]
```

In [54]:
```python
#converting cv data to text
i=0
list_of_sentance_cv=[]
for sentence in X_cv:
    list_of_sentance_cv.append(sentence.split())
```

In [55]:
```python
# 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
```

```python
        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])
```

```
100%|████████████████████████████████████████████████████████████████████
████████| 12041/12041 [00:23<00:00, 519.33it/s]
```

```
(12041, 50)
[ 0.22061323  0.21235475 -0.16896268  0.22343086 -0.09071587  0.1073907
8
 -0.34886903  0.5595728   0.12632388 -0.03554053  0.22296681  0.2554576
1
 -0.63874193 -0.24312705 -0.62764567 -0.66725767  0.41865172  0.3312491
7
  0.0702715  -0.01597985 -0.54764683 -0.27855796  0.03624687 -0.0700586
4
  0.12073137 -0.34791544 -0.1194563   0.34091422 -0.07003629  0.0721434
2
  0.20709204 -0.27818091  0.59850776  0.06985964  0.18257489  0.4814907
5
  0.92962703 -0.43136631 -0.09380837  0.23400903 -0.33236421  0.8882408
4
  0.3867541   0.16926486 -0.13988762 -0.27271803  0.42890321 -0.1445690
4
  0.12084347  0.09244303]
```
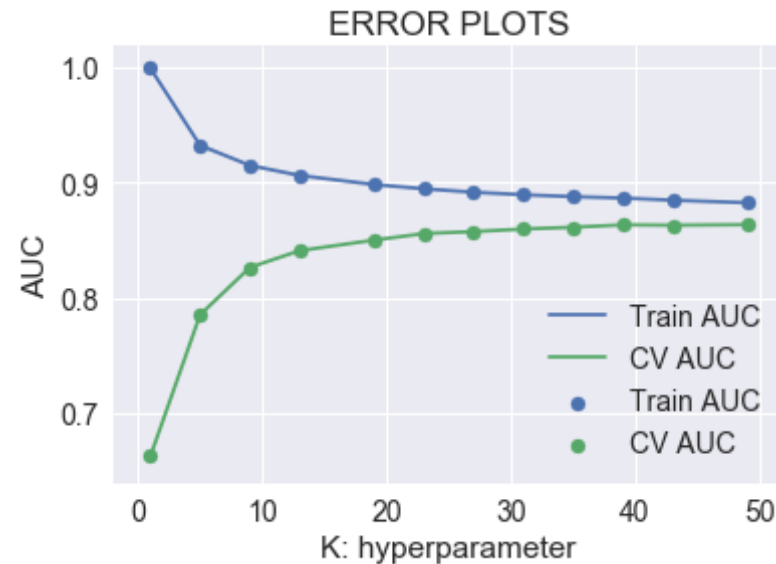
In [56]: 
```python
#Converting Test data text
i=0
```

```
list_of_sentance_test=[]
for sentance in X_test:
    list_of_sentance_test.append(sentance.split())
```

In [57]:
```
# average Word2Vec
# compute average word2vec for each review.
sent_vectors_test = []; # the avg-w2v for each sentence/review is store
d 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, 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])
```

```
100%|████████████████████████████████████████████████████████████
████████| 17971/17971 [15:42<00:00, 19.07it/s]
```

```
(17971, 50)
[ 0.07082332  0.14744843 -0.34326744  0.26433748 -0.25696143  0.2674675
6
 -0.04612008  0.65757128 -0.24657265 -0.00836551  0.29664648  0.2638051
4
 -0.52069131  0.14866833 -0.22348401 -0.27214561  0.67488895  0.4363935
  0.12844289 -0.00646619 -0.1644337   0.1080099  -0.10456489 -0.1174851
  0.33461618 -0.30477659 -0.26067668  0.61385595 -0.17979421  0.1150468
6
 -0.03476828 -0.24744601  0.10591709 -0.1247674   0.32996377  0.5205257
1
  0.65893651 -0.56798559  0.12045094 -0.24029761 -0.74464108  0.2038755
```

```
1
  0.42467368 -0.10131693 -0.16350118  0.16876148  0.39578185 -0.0588788
1
 -0.18479102  0.0384316 ]
```

In [58]:
```python
train_auc = []
cv_auc = []
K = [1,5,9,13,19,23,27,31,35,39,43,49]
for i in K:
    neigh = KNeighborsClassifier(n_neighbors=i,algorithm='brute')
    neigh.fit(sent_vectors_train, y_train)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probab
ility estimates of the positive class
    # not the predicted outputs
    y_train_pred =  neigh.predict_proba(sent_vectors_train)[:,1]
    y_cv_pred =  neigh.predict_proba(sent_vectors_cv)[:,1]

    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(K, train_auc, label='Train AUC')
plt.scatter(K, train_auc, label='Train AUC')
plt.plot(K, cv_auc, label='CV AUC')
plt.scatter(K, cv_auc, label='CV AUC')
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```
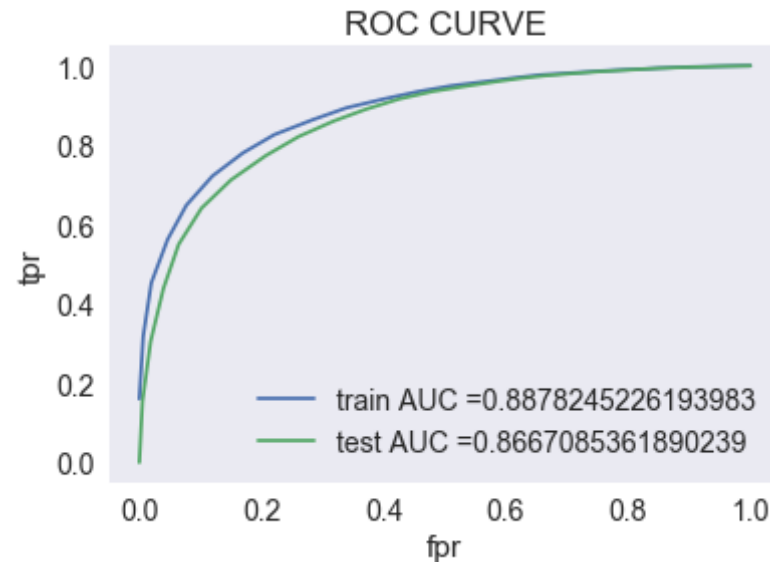
## ERROR PLOTS



```
In [59]: neigh = KNeighborsClassifier(n_neighbors=35,algorithm='brute')
         neigh.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

         y_train_pred = neigh.predict_proba(sent_vectors_train)[:,1]
         y_test_pred = neigh.predict_proba(sent_vectors_test)[:,1]

         train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
         test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

         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("fpr")
         plt.ylabel("tpr")
         plt.title("ROC CURVE")
```

```
plt.grid()
plt.show()
```

## ROC CURVE



```
In [60]: print("="*100)
         from sklearn.metrics import confusion_matrix
         best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
         print("Train confusion matrix")
         print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_
         t)))
         print("Test confusion matrix")
         print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t
         )))
```

```
================================================================================
==============================
the maximum value of tpr*(1-fpr) 0.647890039572986 for threshold 0.829
Train confusion matrix
[[ 3303   673]
 [ 4505 15963]]
Test confusion matrix
[[ 2353   614]
```

```
[ 3401 11603]]
```

In [61]:
```
df_cm = pd.DataFrame((confusion_matrix(y_test, predict_with_best_t(y_te
st_pred, best_t))), range(2),range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
```

Out[61]: `<matplotlib.axes._subplots.AxesSubplot at 0x527a85bb38>`



## [5.1.4] Applying KNN brute force on TFIDF W2V, SET 4

In [62]:
```
# Please write all the code with proper documentation
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
```

```
print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)
```

```
(24444,) (24444,)
(12041,) (12041,)
(17971,) (17971,)
```

In [63]:
```
# Please write all the code with proper documentation
# 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_)))
```

In [64]:
```
# 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_train_sent_vectors = []; # 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_train_sent_vectors.append(sent_vec)
        row += 1
```

```
100%|████████████████████████████████████████████████████████████████████
█████████| 24444/24444 [49:46<00:00,  8.18it/s]
```

In [65]:
```python
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer()
tf_idf_matrix = model.fit_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_)))
```

In [66]:
```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_cv_sent_vectors = []; # 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
            # 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
```
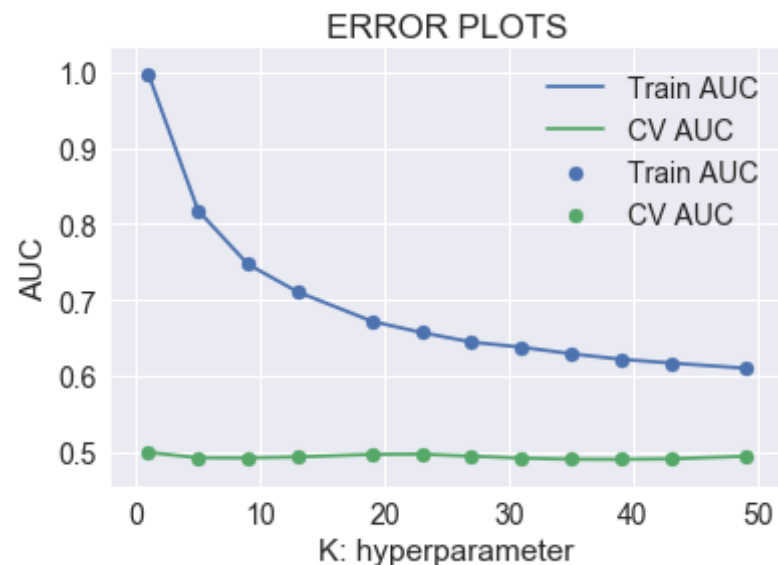
```
            tfidf_cv_sent_vectors.append(sent_vec)
            row += 1
```

```
100%|████████████████████████████████████████████████████████████████████████████| 12041/12041 [02:34<00:00, 78.09it/s]
```

In [67]:
```python
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer()
tf_idf_matrix = model.fit_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_)))
```

In [68]:
```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 cell_val = tfidf

tfidf_test_sent_vectors = []; # 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/review
    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_test_sent_vectors.append(sent_vec)
        row += 1
```

```
100%|████████████████████████████████████████████████████████████
████████████| 17971/17971 [05:10<00:00, 57.93it/s]
```

In [69]:
```python
train_auc = []
cv_auc = []
K = [1,5,9,13,19,23,27,31,35,39,43,49]
for i in K:
    neigh = KNeighborsClassifier(n_neighbors=i,algorithm='brute')
    neigh.fit(tfidf_train_sent_vectors, y_train)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probab
ility estimates of the positive class
    # not the predicted outputs
    y_train_pred =  neigh.predict_proba(tfidf_train_sent_vectors)[:,1]
    y_cv_pred =  neigh.predict_proba(tfidf_cv_sent_vectors)[:,1]

    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(K, train_auc, label='Train AUC')
plt.scatter(K, train_auc, label='Train AUC')
plt.plot(K, cv_auc, label='CV AUC')
plt.scatter(K, cv_auc, label='CV AUC')
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```

ERROR PLOTS

In [70]:
```python
neigh = KNeighborsClassifier(n_neighbors=39,algorithm='brute')
neigh.fit(tfidf_train_sent_vectors, 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

y_train_pred = neigh.predict_proba(tfidf_train_sent_vectors)[:,1]
y_test_pred = neigh.predict_proba(tfidf_test_sent_vectors)[:,1]

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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("fpr")
plt.ylabel("tpr")
plt.title("ROC CURVE")
```

```
plt.grid()
plt.show()
```

ROC CURVE



```
In [71]: print("="*100)
         from sklearn.metrics import confusion_matrix
         best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
         print("Train confusion matrix")
         print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_
         t)))
         print("Test confusion matrix")
         print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t
         )))
```

```
========================================================================
==============================
the maximum value of tpr*(1-fpr) 0.34876833514096617 for threshold 0.84
6
Train confusion matrix
[[ 2387  1637]
 [ 8414 12006]]
Test confusion matrix
[[1292  1641]
```

```
[[1283 1641]
 [6686 8361]]
```

In [72]:
```python
df_cm = pd.DataFrame((confusion_matrix(y_test, predict_with_best_t(y_te
st_pred, best_t))), range(2),range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
```

Out[72]: `<matplotlib.axes._subplots.AxesSubplot at 0x526d7f96d8>`



## [5.2] Applying KNN kd-tree

### [5.2.1] Applying KNN kd-tree on BOW, SET 5

In [73]:
```python
# Please write all the code with proper documentation
labels=final['Score'].values

# Picking the top 3000 points as TSNE
X = final['CleanedText'].values
Y = labels[0:60000]
```

```python
In [74]:  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


          print(X_train.shape, y_train.shape)
          print(X_cv.shape, y_cv.shape)
          print(X_test.shape, y_test.shape)

          from sklearn.feature_extraction.text import CountVectorizer
          vectorizer = CountVectorizer(min_df=10, max_features=500)
          vectorizer.fit(X_train) # fit has to happen only on train data

          # 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)
```

```
(24444,) (24444,)
(12041,) (12041,)
(17971,) (17971,)
After vectorizations
(24444, 500) (24444,)
(12041, 500) (12041,)
(17971, 500) (17971,)
=========================================================================
============================
```
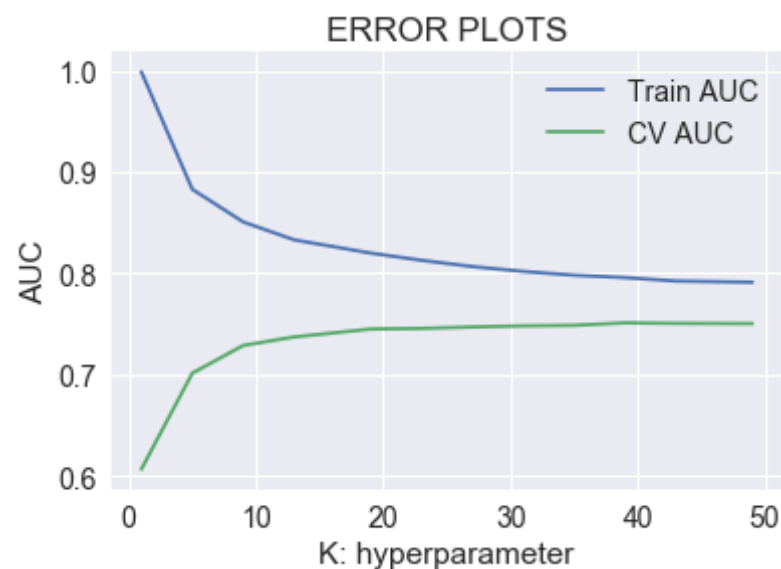
```python
In [75]: train_auc = []
         cv_auc = []
         K = [1,5,9,13,19,23,27,31,35,39,43,49]
         for i in K:
             neigh = KNeighborsClassifier(n_neighbors=i,algorithm='kd_tree')
             neigh.fit(X_train_bow, y_train)
             # roc_auc_score(y_true, y_score) the 2nd parameter should be probab
         ility estimates of the positive class
             # not the predicted outputs
             y_train_pred =  neigh.predict_proba(X_train_bow)[:,1]
             y_cv_pred =  neigh.predict_proba(X_cv_bow)[:,1]

             train_auc.append(roc_auc_score(y_train,y_train_pred))
             cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

         plt.plot(K, train_auc, label='Train AUC')
         plt.plot(K, cv_auc, label='CV AUC')
         plt.legend()
         plt.xlabel("K: hyperparameter")
         plt.ylabel("AUC")
         plt.title("ERROR PLOTS")
         plt.show()
```

```python
In [76]:   #here we are choosing the best_k based on forloop results
           best_k = 35
```

```python
In [77]:   # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc
           _curve.html#sklearn.metrics.roc_curve
           from sklearn.metrics import roc_curve, auc


           neigh = KNeighborsClassifier(n_neighbors=best_k,algorithm='brute')
           neigh.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

           y_train_pred = neigh.predict_proba(X_train_bow)[:,1]
           y_test_pred = neigh.predict_proba(X_test_bow)[:,1]

           train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
           test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

           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("fpr")
           plt.ylabel("tpr")
           plt.title("ROC CURVE")
           plt.grid()
           plt.show()
```
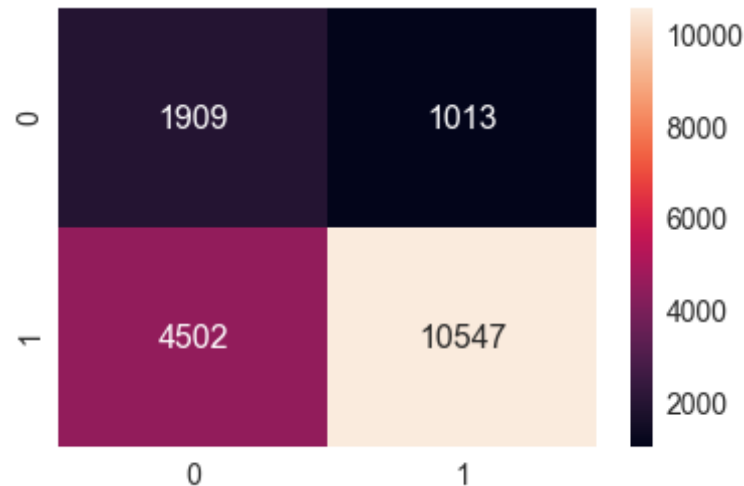
ROC CURVE

```
In [78]:   print("="*100)
           from sklearn.metrics import confusion_matrix
           best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
           print("Train confusion matrix")
           print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_
           t)))
           print("Test confusion matrix")
           print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t
           )))
```

```
========================================================================
==============================
the maximum value of tpr*(1-fpr) 0.5224824871083344 for threshold 0.829
Train confusion matrix
[[ 2892  1062]
 [ 5853 14637]]
Test confusion matrix
[[ 1909  1013]
 [ 4502 10547]]
```

```
In [79]: df_cm = pd.DataFrame((confusion_matrix(y_test, predict_with_best_t(y_te
         st_pred, best_t))), range(2),range(2))
         sns.set(font_scale=1.4)#for label size
         sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
```

```
Out[79]: <matplotlib.axes._subplots.AxesSubplot at 0x526d7fc208>
```



## [5.2.2] Applying KNN kd-tree on TFIDF, <span style="color:red">SET 6</span>

```
In [80]: # Please write all the code with proper documentation
         from sklearn.feature_extraction.text import TfidfVectorizer
         tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
         tf_idf_vect.fit(X_train) # fit has to happen only on train data

         # we use the fitted CountVectorizer to convert the text to vector
         X_train_tfidf = tf_idf_vect.transform(X_train)
         X_cv_tfidf = tf_idf_vect.transform(X_cv)
         X_test_tfidf = tf_idf_vect.transform(X_test)

         print("After vectorizations")
         print(X_train_tfidf.shape, y_train.shape)
         print(X_cv_tfidf.shape, y_cv.shape)
```

```
print(X_test_tfidf.shape, y_test.shape)
print("="*100)
```

After vectorizations
(24444, 14568) (24444,)
(12041, 14568) (12041,)
(17971, 14568) (17971,)
========================================================================
============================

In [81]:
```
train_auc = []
cv_auc = []
K =[1,5,9,13,19,23,27,31,35,39,43,49]
for i in K:
    neigh = KNeighborsClassifier(n_neighbors=i,algorithm='kd_tree')
    neigh.fit(X_train_tfidf, y_train)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probab
ility estimates of the positive class
    # not the predicted outputs
    y_train_pred =  neigh.predict_proba(X_train_tfidf)[:,1]
    y_cv_pred =  neigh.predict_proba(X_cv_tfidf)[:,1]

    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(K, train_auc, label='Train AUC')
plt.plot(K, cv_auc, label='CV AUC')
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```

## ERROR PLOTS



In [82]: 
```python
#here we are choosing the best_k based on forloop results
best_k = 29
```

In [83]: 
```python
from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k,algorithm='kd_tree')
neigh.fit(X_train_tfidf, 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

y_train_pred = neigh.predict_proba(X_train_tfidf)[:,1]
y_test_pred = neigh.predict_proba(X_test_tfidf)[:,1]

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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("fpr")
plt.ylabel("tpr")
plt.title("ROC CURVE")
plt.grid()
plt.show()
```



ROC CURVE

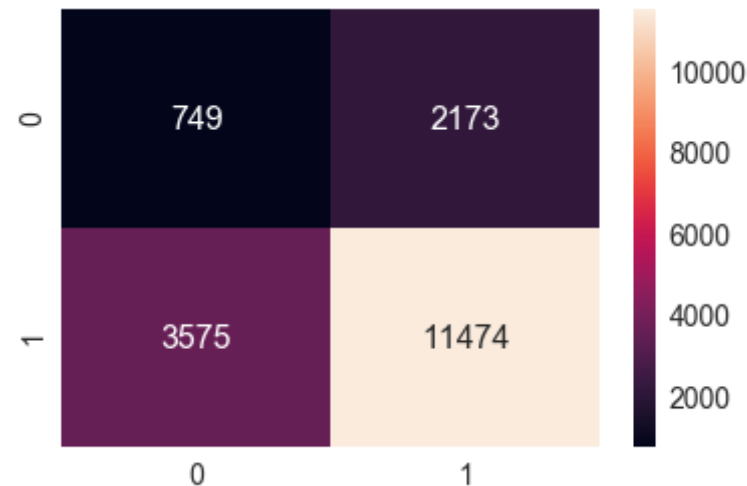train AUC =0.7239310514054624
test AUC =0.5167275606052398

In [84]:
```
print("="*100)
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
```

```
========================================================================
==============================
the maximum value of tpr*(1-fpr) 0.4182881689946833 for threshold 0.897
Train confusion matrix
[[ 2115  1839]
```

```
  [ 4467 16023]]
Test confusion matrix
[[  749  2173]
 [ 3575 11474]]
```

In [85]:
```python
df_cm = pd.DataFrame((confusion_matrix(y_test, predict_with_best_t(y_te
st_pred, best_t))), range(2),range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
```

Out[85]: `<matplotlib.axes._subplots.AxesSubplot at 0x527aafc0f0>`



### [5.2.3] Applying KNN kd-tree on AVG W2V, SET 3

In [86]:
```python
train_auc = []
cv_auc = []
K = [1,5,9,13,19,23,27,31,35,39,43,49]
for i in K:
    neigh = KNeighborsClassifier(n_neighbors=i,algorithm='kd_tree')
    neigh.fit(sent_vectors_train, y_train)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probab
ility estimates of the positive class
```
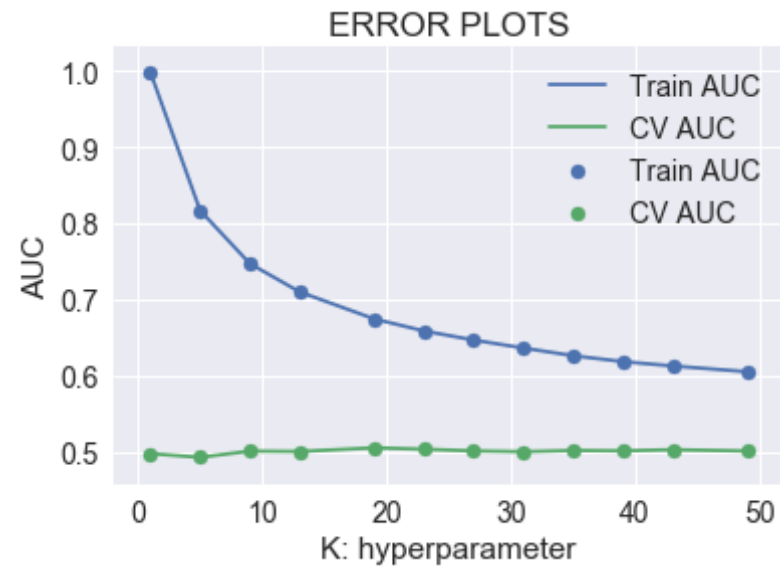
```
    # not the predicted outputs
    y_train_pred =  neigh.predict_proba(sent_vectors_train)[:,1]
    y_cv_pred =  neigh.predict_proba(sent_vectors_cv)[:,1]

    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(K, train_auc, label='Train AUC')
plt.scatter(K, train_auc, label='Train AUC')
plt.plot(K, cv_auc, label='CV AUC')
plt.scatter(K, cv_auc, label='CV AUC')
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```



```
In [87]:  neigh = KNeighborsClassifier(n_neighbors=37,algorithm='kd_tree')
          neigh.fit(sent_vectors_train, y_train)

          y_train_pred = neigh.predict_proba(sent_vectors_train)[:,1]
```
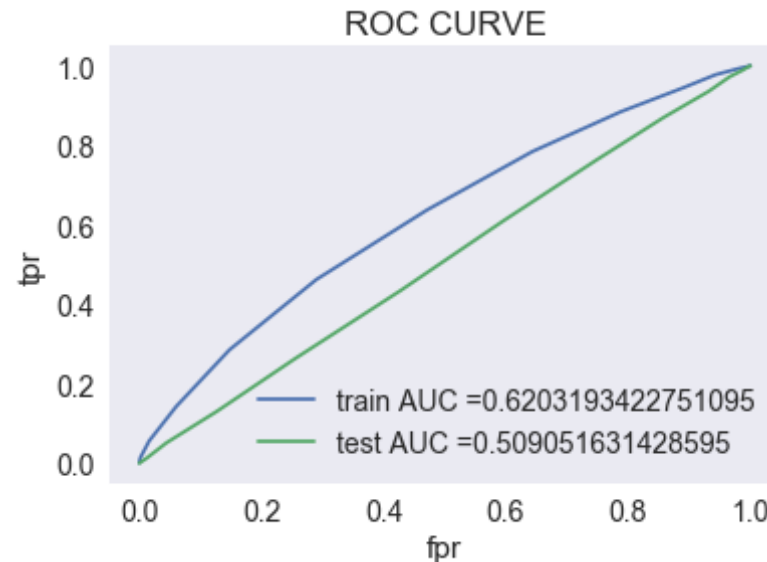
```python
y_test_pred = neigh.predict_proba(sent_vectors_test)[:,1]

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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("fpr")
plt.ylabel("tpr")
plt.title("ROC CURVE")
plt.grid()
plt.show()
```



```python
In [88]: print("="*100)
         from sklearn.metrics import confusion_matrix
         best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
         print("Train confusion matrix")
         print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_
```
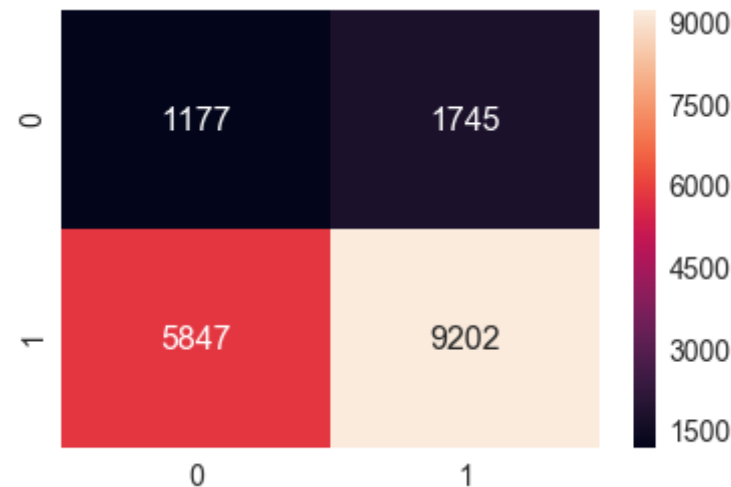
```
t)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t
)))
```

```
================================================================================
==============================
the maximum value of tpr*(1-fpr) 0.3372878883144448 for threshold 0.838
Train confusion matrix
[[ 2093  1861]
 [ 7434 13056]]
Test confusion matrix
[[1177 1745]
 [5847 9202]]
```

In [89]:
```
df_cm = pd.DataFrame((confusion_matrix(y_test, predict_with_best_t(y_te
st_pred, best_t))), range(2),range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
```
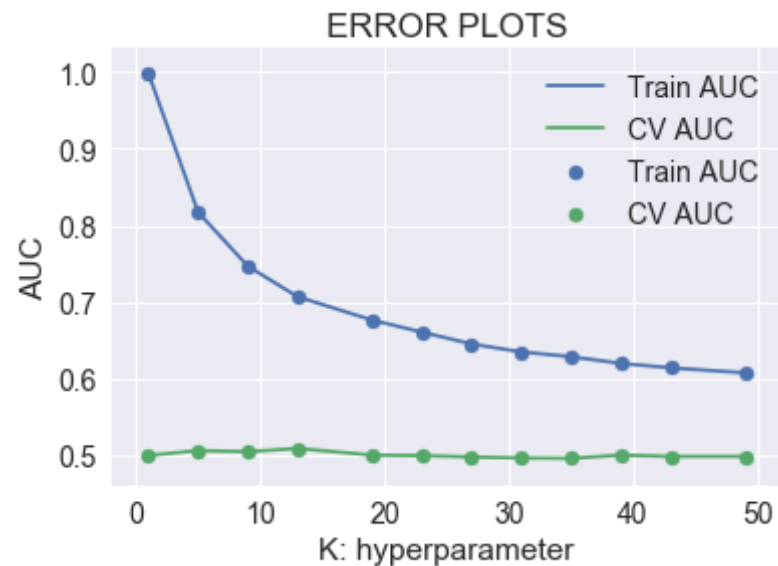
Out[89]: <matplotlib.axes._subplots.AxesSubplot at 0x525248fcc0>



**[5.2.4] Applying KNN kd-tree on TFIDF W2V, SET 4**

```python
In [90]: train_auc = []
         cv_auc = []
         K = [1,5,9,13,19,23,27,31,35,39,43,49]
         for i in K:
             neigh = KNeighborsClassifier(n_neighbors=i,algorithm='kd_tree')
             neigh.fit(tfidf_train_sent_vectors, y_train)
             # roc_auc_score(y_true, y_score) the 2nd parameter should be probab
         ility estimates of the positive class
             # not the predicted outputs
             y_train_pred =  neigh.predict_proba(tfidf_train_sent_vectors)[:,1]
             y_cv_pred =  neigh.predict_proba(tfidf_cv_sent_vectors)[:,1]

             train_auc.append(roc_auc_score(y_train,y_train_pred))
             cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

         plt.plot(K, train_auc, label='Train AUC')
         plt.scatter(K, train_auc, label='Train AUC')
         plt.plot(K, cv_auc, label='CV AUC')
         plt.scatter(K, cv_auc, label='CV AUC')
         plt.legend()
         plt.xlabel("K: hyperparameter")
         plt.ylabel("AUC")
         plt.title("ERROR PLOTS")
         plt.show()
```

ERROR PLOTS



In [95]:
```python
neigh = KNeighborsClassifier(n_neighbors=37,algorithm='kd_tree')
neigh.fit(tfidf_train_sent_vectors, 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

y_train_pred = neigh.predict_proba(tfidf_train_sent_vectors)[:,1]
y_test_pred = neigh.predict_proba(tfidf_test_sent_vectors)[:,1]

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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("fpr")
plt.ylabel("tpr")
plt.title("ROC CURVE")
```
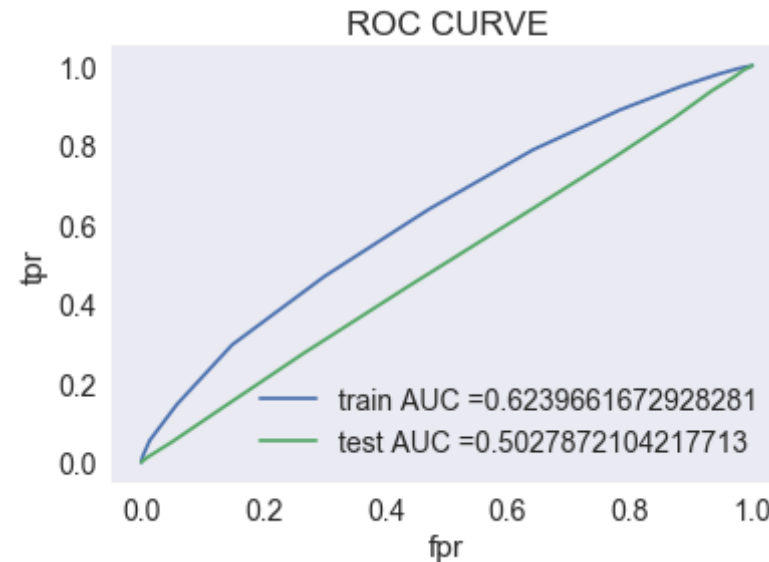
```
plt.grid()
plt.show()
```
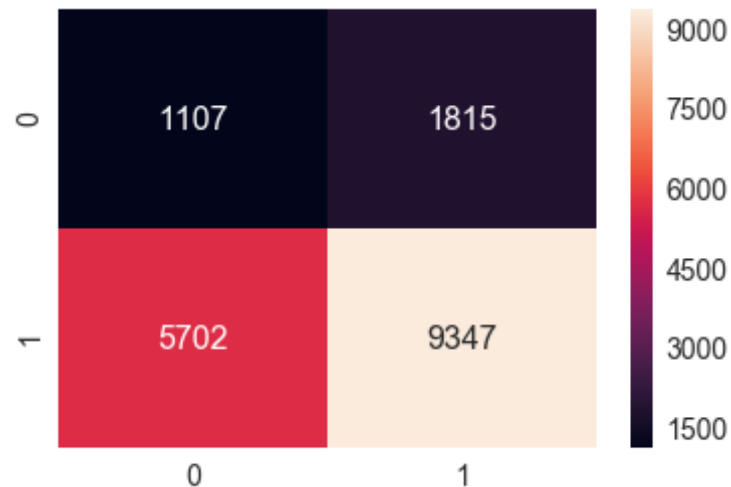
### ROC CURVE



```
In [96]: print("="*100)
         from sklearn.metrics import confusion_matrix
         best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
         print("Train confusion matrix")
         print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_
         t)))
         print("Test confusion matrix")
         print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t
         )))
```

```
====================================================================================
=============================
the maximum value of tpr*(1-fpr) 0.3378562349399747 for threshold 0.838
Train confusion matrix
[[ 2093  1861]
 [ 7412 13078]]
Test confusion matrix
[[1107 1815]
 [5702 0047]]
```

```
[5702 9347]]
```

In [97]:
```python
df_cm = pd.DataFrame((confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))), range(2),range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
```

Out[97]: &lt;matplotlib.axes._subplots.AxesSubplot at 0x52525e12e8&gt;



## [6] Conclusions

In [99]:
```python
# Please compare all your models using Prettytable library
from prettytable import PrettyTable

x = PrettyTable()

x.field_names = ["Vectorizer", "Model", "Hyper parameter", "Train AUC"]

x.add_row(["BOW", "brute", 31,0.77 ])
x.add_row(["TFIDF", "brute", 25, 0.71])
x.add_row(["Avg W2V", "brute", 35, 0.89])
x.add_row(["TFIDF W2V", "brute", 39, 0.62])
```

```python
x.add_row(["BOW", "kd_tree", 35, 0.79])
x.add_row(["TFIDF", "kd_tree", 29, 0.72])
x.add_row(["Avg W2V", "kd_tree", 37, 0.62])
x.add_row(["TFIDF W2V", "kd_tree", 37, 0.62])


print(x)
```

```
+-----------+---------+-----------------+-----------+
| Vectorizer |  Model  | Hyper parameter | Train AUC |
+-----------+---------+-----------------+-----------+
|    BOW    |  brute  |       31        |    0.77   |
|   TFIDF   |  brute  |       25        |    0.71   |
|  Avg W2V  |  brute  |       35        |    0.89   |
| TFIDF W2V |  brute  |       39        |    0.62   |
|    BOW    | kd_tree |       35        |    0.79   |
|   TFIDF   | kd_tree |       29        |    0.72   |
|  Avg W2V  | kd_tree |       37        |    0.62   |
| TFIDF W2V | kd_tree |       37        |    0.62   |
+-----------+---------+-----------------+-----------+
```