

Amazon Fine Food Reviews Analysis

Data Source: <https://www.kaggle.com/snap/amazon-fine-food-reviews> (<https://www.kaggle.com/snap/amazon-fine-food-reviews>)

EDA: <https://nycdatasience.com/blog/student-works/amazon-fine-foods-visualization/>
(<https://nycdatasience.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:

ProductId - unique identifier for the product UserId - unique identifier for the user ProfileName
HelpfulnessNumerator - number of users who found the review helpful HelpfulnessDenominator - number of users who indicated whether they found the review helpful or not Score - rating between 1 and 5 Time - timestamp for the review Summary - brief summary of the review 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 the Score/Rating. A rating of 4 or 5 could be considered a positive review. A review of 1 or 2 could be considered negative. A review of 3 is neutral and ignored. This is an approximate and proxy way of determining the polarity (positivity/negativity) of a review.

Loading the data

The dataset is available in two forms

.csv file 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 will 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
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

from sklearn.model_selection import train_test_split
from sklearn.metrics import roc_auc_score
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import accuracy_score
from sklearn.cross_validation import cross_val_score
from collections import Counter
from sklearn import cross_validation

```

D:\Anaconda3\lib\site-packages\gensim\utils.py:1209: UserWarning: detected Windows; aliasing chunkize to chunkize_serial

warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")

D:\Anaconda3\lib\site-packages\sklearn\cross_validation.py:41: DeprecationWarning: This module was deprecated in version 0.18 in favor of the model_selection module into which all the refactored classes and functions are moved. Also note that the interface of the new CV iterators are different from that of this module. This module will be removed in 0.20.

"This module will be removed in 0.20.", DeprecationWarning)

[1]. Reading Data

```
In [2]: # using SQLite Table to read data.
con = sqlite3.connect('D:\\TGM\\ML\\AmazonFineFoodReviews\\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 500000 data points
# you can change the number to any other number based on your computing power

# filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3
LIMIT 500000""", con)
# for tsne assignment you can take 5k data points

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

# Give reviews with Score>3 a positive rating(1), and reviews with a score<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 (525814, 10)

Out[2]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	Helpfulne
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0
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	Cou
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 [5]: display[display['UserId']=='AZY10LLTJ71NX']
```

```
Out[5]:
```

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

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

```
Out[6]: 393063
```

Exploratory Data Analysis

[2] 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	Helpful
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 can be seen above the same user has multiple reviews of the with the 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 delete 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=True, 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]: (364173, 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]: 69.25890143662969
```

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 calculations

```
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	Helpfulr
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 entrie
s left
print(final.shape)

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

(364171, 10)
```

```
Out[13]: 1    307061
0     57110
Name: Score, dtype: int64
```

[3]. Text Preprocessing.

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 observed to be better than Porter Stemming)

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

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

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

this witty little book makes my son laugh at loud. i recite it in the car as we're driving along and he always can sing the refrain. he's learned about wh ales, India, drooping roses: i love all the new words this book introduces and the silliness of it all. this is a classic book i am willing to bet my son will STILL be able to recite from memory when he is in college

=====

I was really looking forward to these pods based on the reviews. Starbucks is good, but I prefer bolder taste.... imagine my surprise when I ordered 2 boxes - both were expired! One expired back in 2005 for gosh sakes. I admit that Amazon agreed to credit me for cost plus part of shipping, but geez, 2 years expired!!! I'm hoping to find local San Diego area shoppe that carries pods so that I can try something different than starbucks.

=====

Great ingredients although, chicken should have been 1st rather than chicken broth, the only thing I do not think belongs in it is Canola oil. Canola or rapeseed is not something a dog would ever find in nature and if it did find rapeseed in nature and eat it, it would poison them. Today's Food industries have convinced the masses that Canola oil is a safe and even better oil than olive or virgin coconut, facts though say otherwise. Until the late 70's it was poisonous until they figured out a way to fix that. I still like it but it could be better.

=====

Can't do sugar. Have tried scores of SF Syrups. NONE of them can touch the excellence of this product.

Thick, delicious. Perfect. 3 ingredients: Water, Maltitol, Natural Maple Flavor. PERIOD. No chemicals. No garbage.

Have numerous friends & family members hooked on this stuff. My husband & son, who do NOT like "sugar free" prefer this over major label regular syrup.

I use this as my SWEETENER in baking: cheesecakes, white brownies, muffins, pumpkin pies, etc... Unbelievably delicious...

Can you tell I like it? :)

=====

```
In [15]: # remove urls from text python: https://stackoverflow.com/a/40823105/4084039
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)
```

this witty little book makes my son laugh at loud. i recite it in the car as we're driving along and he always can sing the refrain. he's learned about whales, India, drooping roses: i love all the new words this book introduces and the silliness of it all. this is a classic book i am willing to bet my son will STILL be able to recite from memory when he is in college

```
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)
```

this witty little book makes my son laugh at loud. i recite it in the car as we're driving along and he always can sing the refrain. he's learned about wh ales, India, drooping roses: i love all the new words this book introduces and the silliness of it all. this is a classic book i am willing to bet my son will STILL be able to recite from memory when he is in college

=====

I was really looking forward to these pods based on the reviews. Starbucks is good, but I prefer bolder taste.... imagine my surprise when I ordered 2 boxes - both were expired! One expired back in 2005 for gosh sakes. I admit that Amazon agreed to credit me for cost plus part of shipping, but geez, 2 years expired!!! I'm hoping to find local San Diego area shoppe that carries pods so that I can try something different than starbucks.

=====

Great ingredients although, chicken should have been 1st rather than chicken broth, the only thing I do not think belongs in it is Canola oil. Canola or rapeseed is not something a dog would ever find in nature and if it did find rapeseed in nature and eat it, it would poison them. Today's Food industries have convinced the masses that Canola oil is a safe and even better oil than olive or virgin coconut, facts though say otherwise. Until the late 70's it was poisonous until they figured out a way to fix that. I still like it but it could be better.

=====

Can't do sugar. Have tried scores of SF Syrups. NONE of them can touch the excellence of this product. Thick, delicious. Perfect. 3 ingredients: Water, Maltitol, Natural Maple Flavor. PERIOD. No chemicals. No garbage. Have numerous friends & family members hooked on this stuff. My husband & son, who do NOT like "sugar free" prefer this over major label regular syrup. I use this as my SWEETENER in baking: cheesecakes, white brownies, muffins, pumpkin pies, etc... Unbelievably delicious...Can you tell I like it? :)

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

Great ingredients although, chicken should have been 1st rather than chicken broth, the only thing I do not think belongs in it is Canola oil. Canola or rapeseed is not something a dog would ever find in nature and if it did find rapeseed in nature and eat it, it would poison them. Today's food industries have convinced the masses that Canola oil is a safe and even better oil than olive or virgin coconut, facts though say otherwise. Until the late 70s it was as poisonous until they figured out a way to fix that. I still like it but it could be better.

=====

In [19]: *#remove words with numbers python: <https://stackoverflow.com/a/18082370/4084039>*
`sent_0 = re.sub("\S*\d\S*", "", sent_0).strip()`
`print(sent_0)`

this witty little book makes my son laugh at loud. i recite it in the car as we're driving along and he always can sing the refrain. he's learned about whales, India, drooping roses: i love all the new words this book introduces and the silliness of it all. this is a classic book i am willing to bet my son will STILL be able to recite from memory when he is in college

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

Great ingredients although chicken should have been 1st rather than chicken broth the only thing I do not think belongs in it is Canola oil Canola or rape seed is not something a dog would ever find in nature and if it did find rape seed in nature and eat it it would poison them Today's food industries have convinced the masses that Canola oil is a safe and even better oil than olive or virgin coconut facts though say otherwise Until the late 70s it was poisonous until they figured out a way to fix that I still like it but it could be better

```
In [21]: # 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 been removed 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', "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', "isn'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]: # Combining all the above students
from tqdm import tqdm
preprocessed_reviews = []
# tqdm is for printing the status bar
for sentence in tqdm(final['Text'].values):
    sentence = re.sub(r"http\S+", "", sentence)
    sentence = BeautifulSoup(sentence, 'lxml').get_text()
    sentence = decontracted(sentence)
    sentence = re.sub("\S*\d\S*", "", sentence).strip()
    sentence = re.sub('[^A-Za-z]+', ' ', sentence)
    # https://gist.github.com/sebleier/554280
    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not in stopwords)
    preprocessed_reviews.append(sentence.strip())
```

```
100%|██████████| 364171/364171 [06:25<00:00, 944.81it/s]
```

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

```
Out[23]: 'great ingredients although chicken rather chicken broth thing not think belo
ngs canola oil canola rapeseed not someting dog would ever find nature find r
apeseed nature eat would poison today food industries convinced masses canola
oil safe even better oil olive virgin coconut facts though say otherwise late
poisonous figured way fix still like could better'
```

```
In [24]: final['cleaned text']=preprocessed reviews
```

```
In [25]: final.shape
```

```
Out[25]: (364171, 11)
```

```
In [26]: data_pos = final[final["Score"] == 1].sample(n = 10000)
data_neg = final[final["Score"] == 0].sample(n = 10000)
final1 = pd.concat([data_pos, data_neg])
final1.shape
```

```
Out[26]: (20000, 11)
```

```
In [27]: Y = final1['Score'].values
X = final1['cleaned_text'].values
print(Y.shape)
print(type(Y))
print(X.shape)
print(type(X))
```

```
(20000,)
<class 'numpy.ndarray'>
(20000,)
<class 'numpy.ndarray'>
```

```

In [28]: # split the data set into train and test
X_Train, X_Test, Y_Train, Y_Test = train_test_split(X,Y,test_size=0.3, random_
state=0)

# split the train data set into cross validation train and cross validation te
st
X_tr, X_cv, Y_tr, Y_cv = train_test_split(X,Y, test_size=0.3, random_state=0)

#X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.2,random_state
=12)
#X_train,X_cv,Y_train,Y_cv=train_test_split(X,Y,test_size=0.2,random_state=12)
print('='*100)
print("After splitting")
print("X_Train Shape:",X_Train.shape,"Y_Train Shape:",Y_Train.shape)
print("X_cv Shape:",X_cv.shape,"Y_cv Shape",Y_cv.shape)
print("X_Test Shape",X_Test.shape,"Y_Test Shape",Y_Test.shape)

=====
=====
After splitting
X_Train Shape: (14000,) Y_Train Shape: (14000,)
X_cv Shape: (6000,) Y_cv Shape (6000,)
X_Test Shape (6000,) Y_Test Shape (6000,)

```

```

In [29]: # creating odd list of K for KNN
#myList = list(range(80,100,2))
#neighbors = list(filter(lambda x: x % 2 != 0, myList))
def Optimal_K(X_Train, Y_Train):
    neighbors = list(range(1,50,2))
    # empty list that will hold cv scores
    cv_scores = []

    # perform 10-fold cross validation
    for k in neighbors:
        knn = KNeighborsClassifier(n_neighbors=k, algorithm = 'brute')
        scores = cross_val_score(knn, X_Train, Y_Train, cv=10, scoring='ac
curacy')
        cv_scores.append(scores.mean())

    # changing to misclassification error
    MSE = [1 - x for x in cv_scores]

    # determining best k
    bestK = neighbors[MSE.index(min(MSE))]
    print('\nThe optimal number of neighbors is %d.' % bestK)

    # plot misclassification error vs k
    plt.plot(neighbors, MSE)

    for xy in zip(neighbors, np.round(MSE,3)):
        plt.annotate('%s, %s' % xy, xy=xy, textcoords='data')

    plt.xlabel('Number of Neighbors K')
    plt.ylabel('Misclassification Error')
    plt.show()

    print("the misclassification error for each k value is : ", np.round(M
SE,3))
    return bestK

```

[4] Featurization - Using Brute Force Algorithm

[4.1] Bag Of Words

```

In [30]: #Bow
count_vect = CountVectorizer() #in scikit-learn
count_vect.fit(X_Train)
print("some feature names ", count_vect.get_feature_names()[:10])
X_Train_Bow = count_vect.transform(X_Train)
X_Test_Bow = count_vect.transform(X_Test)
X_CV_Bow = count_vect.transform(X_cv)

print('='*50)

#final_counts = count_vect.transform(X_Test)

print("the type of X Train : ",type(X_Train_Bow))
print("the shape of Train BOW vectorizer ",X_Train_Bow.get_shape())
print("the shape of Test BOW vectorizer ",X_Test_Bow.get_shape())
print("the shape of CV BOW vectorizer ",X_CV_Bow.get_shape())
#print("the number of unique words ", final_counts.get_shape()[1])

some feature names  ['aa', 'ab', 'aback', 'abandon', 'abandoned', 'abandonin
g', 'abba', 'abby', 'abbypomeroy', 'abc']
=====
the type of X Train :  <class 'scipy.sparse.csr.csr_matrix'>
the shape of Train BOW vectorizer  (14000, 23954)
the shape of Test BOW vectorizer  (6000, 23954)
the shape of CV BOW vectorizer  (6000, 23954)

```

[4.1.1] AUC Curve Plot

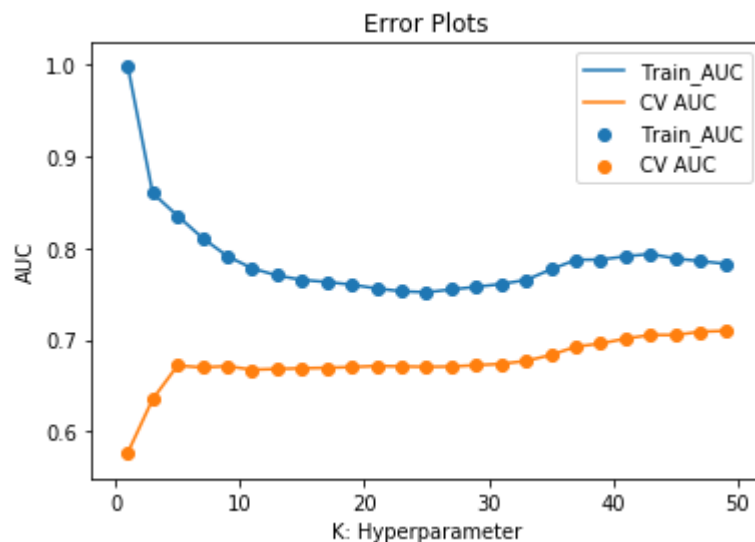
```

In [31]: train_AUC = []
        CV_AUC = []

        K = list(range(1,50,2))
        for i in K:
            knn = KNeighborsClassifier(n_neighbors = i, algorithm = 'brute')
            #fit a model on train BOW vectorizer
            knn.fit(X_Train_Bow, Y_Train)
            #predict probabilities on train BOW vectorizer
            Y_Train_Pred = knn.predict_proba(X_Train_Bow)[: ,1]
            #predict probabilities on Cross validation BOW vectorizer
            Y_CV_Pred = knn.predict_proba(X_CV_Bow)[: ,1]
            #calculate AUC score
            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()

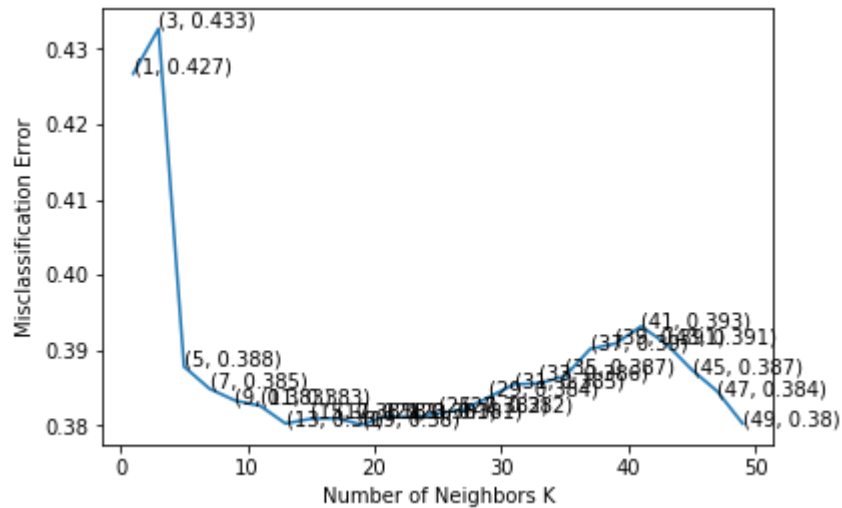
```



[4.1.2] 10-fold cross validation, determining best K

```
In [32]: optimal_k1 = Optimal_K(X_Train_Bow, Y_Train)
print("optimal_k1:", optimal_k1)
```

The optimal number of neighbors is 19.



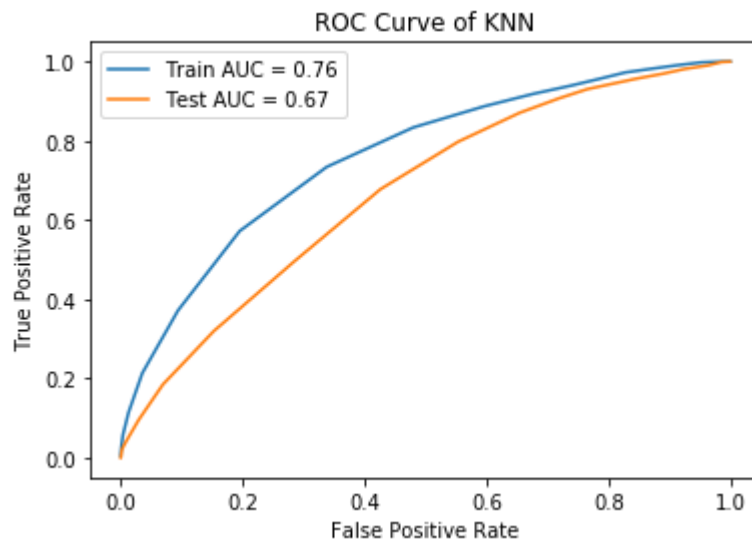
the misclassification error for each k value is : [0.427 0.433 0.388 0.385
0.383 0.383 0.38 0.381 0.381 0.38 0.381 0.381
0.382 0.382 0.384 0.385 0.386 0.387 0.39 0.391 0.393 0.391 0.387 0.384
0.38]
optimal_k1: 19

```
In [33]: optimal_model = KNeighborsClassifier(n_neighbors=optimal_k1, algorithm='brute'
)
optimal_model.fit(X_Train_Bow, Y_Train)
prediction = optimal_model.predict(X_Test_Bow)
```

[4.1.3] ROC Curve of KNN

```
In [34]: #with the reference of below link:
#https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn
#machine-learning-algorithm-using-python-and-sci
#predict probabilities on X_Train_Bow and X_Test_Bow and pass as param to roc_
#curve to find roc curve
Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, optimal_model.predict_proba(X_Train_Bow)[:,-1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, optimal_model.predict_proba(X_Test_Bow)[:,-1])
roc_auc = auc(Train_FPR, Train_TPR)
roc_auc1 = auc(Test_FPR, Test_TPR)

plt.plot(Train_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc1)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of KNN')
plt.show()
```



[4.1.4]Train and Test Accuracy

```
In [35]: Training_Accuracy = optimal_model.score(X_Train_Bow, Y_Train)
print('Training_Accuracy=%0.3f'%Training_Accuracy)
Training_Error = 1 - Training_Accuracy
print('Training_Error=%0.3f'%Training_Error)

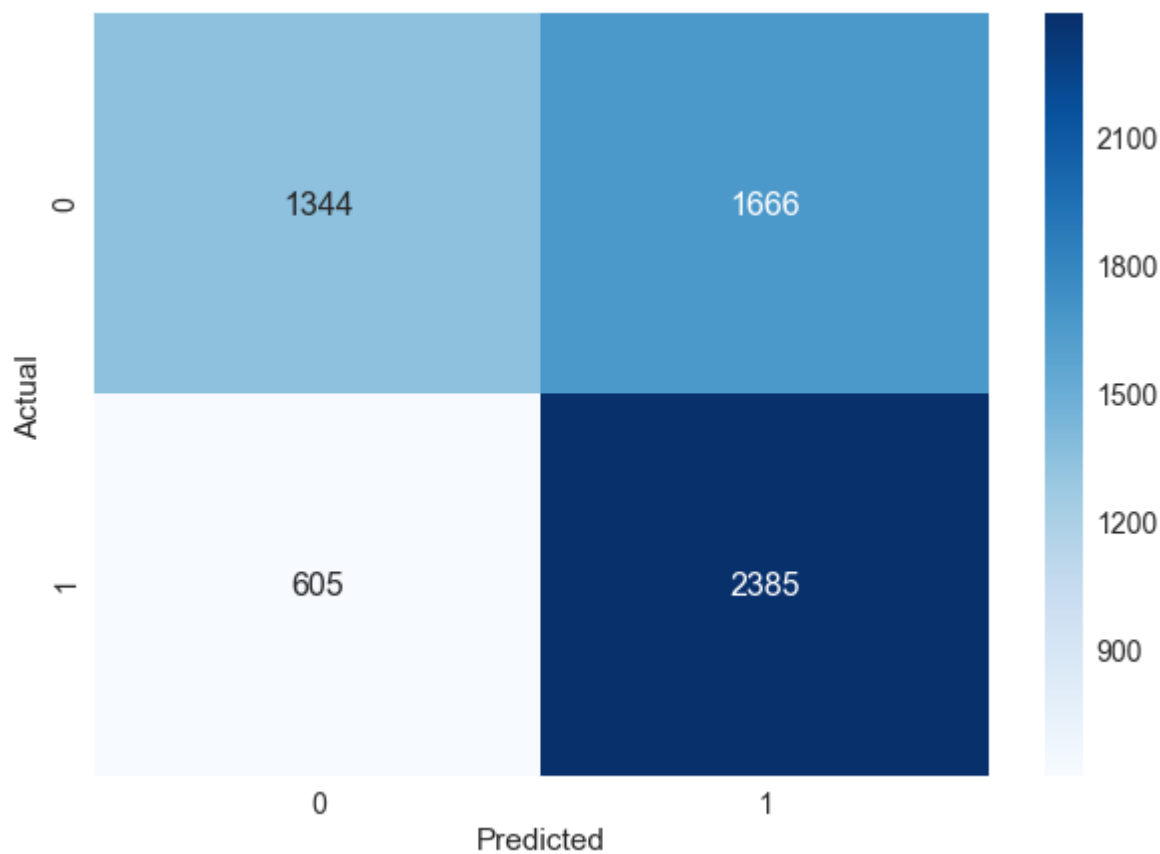
Test_Accuracy = accuracy_score(Y_Test, prediction)
print('Test_Accuracy=%0.3f'%Test_Accuracy)
Test_Error = 1 - Test_Accuracy
print('Test_Error=%0.3f'%Test_Error)

Training_Accuracy=0.677
Training_Error=0.323
Test_Accuracy=0.622
Test_Error=0.378
```

[4.1.5] Confusion Matrix

```
In [36]: #With the reference of below link:
#https://www.kaggle.com/agungor2/various-confusion-matrix-plots
from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, optimal_model.predict(X_Test_Bow))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np
.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, f
mt='d')
```

Out[36]: <matplotlib.axes._subplots.AxesSubplot at 0x2268082d358>



[4.1.6] Classification Report


```
In [37]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction))
```

	precision	recall	f1-score	support
0	0.69	0.45	0.54	3010
1	0.59	0.80	0.68	2990
avg / total	0.64	0.62	0.61	6000

[4.2] TF-IDF

```
In [38]: #TF-IDF
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=5)
tf_idf_vect.fit_transform(X_Train)
print("some sample features(unique words in the corpus)",tf_idf_vect.get_featu
re_names()[0:10])
print('='*50)

X_Train_TfIdf = count_vect.transform(X_Train)
X_Test_TfIdf = count_vect.transform(X_Test)
X_CV_TfIdf = count_vect.transform(X_cv)

#final_tf_idf = tf_idf_vect.transform(X_Test)
print("the type of count vectorizer ",type(X_Train_TfIdf))
print("the shape of out text TFIDF vectorizer ",X_Train_TfIdf.get_shape())
print("the shape of out text TFIDF vectorizer ",X_Test_TfIdf.get_shape())
print("the shape of out text TFIDF vectorizer ",X_CV_TfIdf.get_shape())
#print("the number of unique words including both unigrams and bigrams ", fina
l_tf_idf.get_shape()[1])
```

```
some sample features(unique words in the corpus) ['aa', 'abdominal', 'abilit
y', 'able', 'able buy', 'able drink', 'able eat', 'able enjoy', 'able find',
'able finish']
```

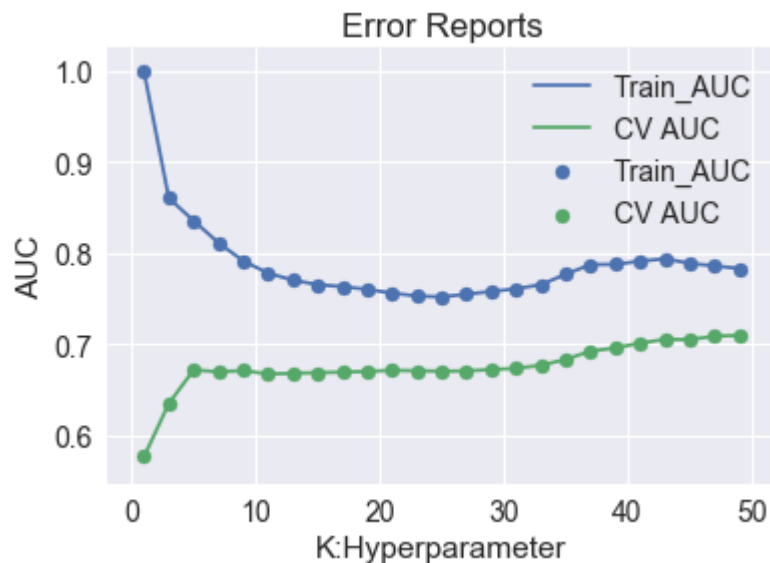
```
=====
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer (14000, 23954)
the shape of out text TFIDF vectorizer (6000, 23954)
the shape of out text TFIDF vectorizer (6000, 23954)
```

[4.2.1] AUC Curve Plot

```

In [39]: train_AUC=[]
CV_AUC=[]
K = list(range(1,50,2))
for i in K:
    knn = KNeighborsClassifier(n_neighbors=i, algorithm='brute')
    #fit a model on Train TFIDF vectorizer
    knn.fit(X_Train_TfIdf, Y_Train)
    #predict probabilities on train TfIdf vectorizer
    Y_Train_Pred = knn.predict_proba(X_Train_TfIdf)[:,:1]
    #predict probabilities on cross validation TfIdf vectorizer
    Y_CV_Pred = knn.predict_proba(X_CV_TfIdf)[:,:1]
    #calculate AUC score
    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 Reports')
plt.show()

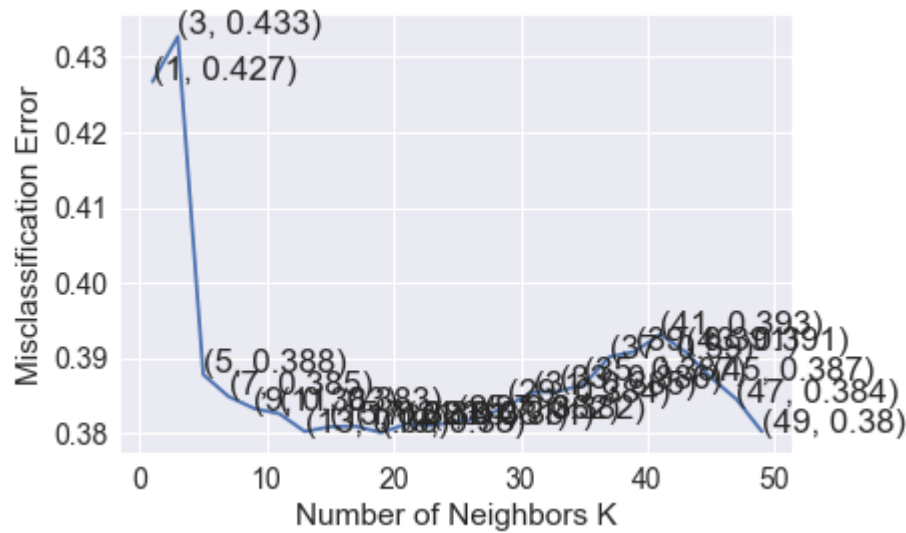
```



[4.2.2] 10-fold cross validation

```
In [40]: optimal_k2 = Optimal_K(X_Train_TfIdf, Y_Train)
print("optimal_k2:", optimal_k2)
```

The optimal number of neighbors is 19.



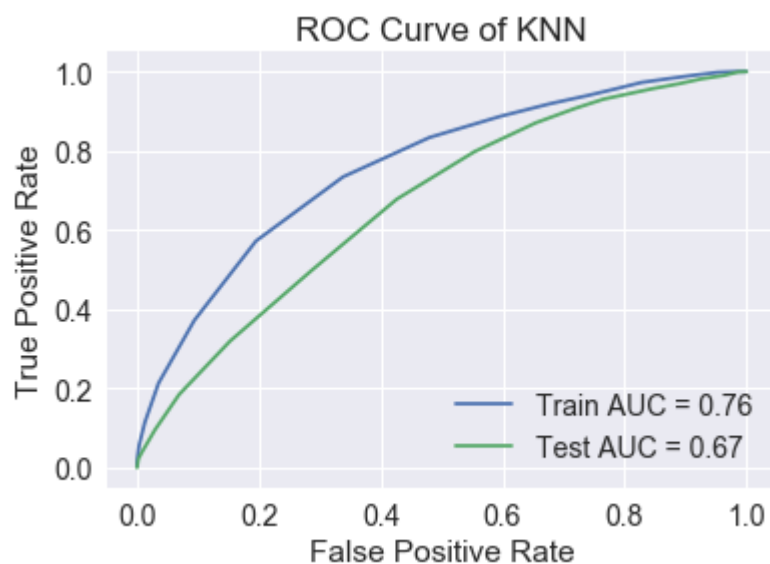
```
the misclassification error for each k value is : [0.427 0.433 0.388 0.385
0.383 0.383 0.38 0.381 0.381 0.38 0.381 0.381
0.382 0.382 0.384 0.385 0.386 0.387 0.39 0.391 0.393 0.391 0.387 0.384
0.38 ]
optimal_k2: 19
```

```
In [41]: optimal_model = KNeighborsClassifier(n_neighbors=optimal_k2, algorithm='brute'
)
optimal_model.fit(X_Train_TfIdf, Y_Train)
prediction = optimal_model.predict(X_Test_TfIdf)
```

[4.2.3] ROC Curve of KNN

```
In [42]: #with the reference of below link:
#https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn
#machine-learning-algorithm-using-python-and-sci
#predict probabilities on X_Train_Bow and X_Test_Bow and pass as param to roc_
#curve to find roc curve
Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, optimal_model.predict_proba(X_Train_Tfidf)[:,:])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, optimal_model.predict_proba(X_Test_Tfidf)[:,:])
roc_auc2 = auc(Train_FPR, Train_TPR)
roc_auc3 = auc(Test_FPR, Test_TPR)

plt.title('Receiving Operating Characteristic Curve')
plt.plot(Train_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc2)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc3)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of KNN')
plt.show()
```



 [4.2.4] Train and Test Accuracy

```
In [43]: Training_Accuracy = optimal_model.score(X_Train_Tfidf, Y_Train)
print('Training_Accuracy=%0.3f'%Training_Accuracy)
Training_Error = 1 - Training_Accuracy
print('Training_Error=%0.3f'%Training_Error)

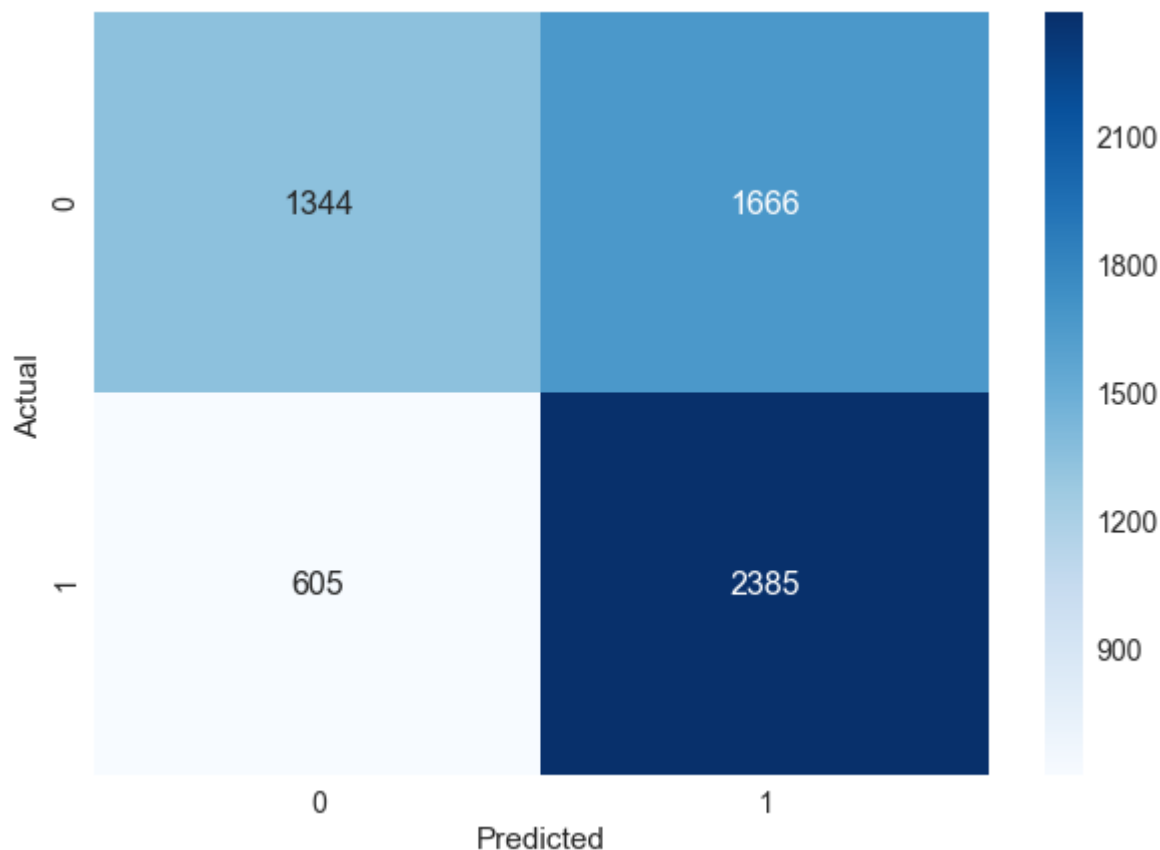
Test_Accuracy = accuracy_score(Y_Test, prediction)
print('Test_Accuracy=%0.3f'%Test_Accuracy)
Test_Error = 1 - Test_Accuracy
print('Test_Error=%0.3f'%Test_Error)

Training_Accuracy=0.677
Training_Error=0.323
Test_Accuracy=0.622
Test_Error=0.378
```

[4.2.5]Confusion Matrix

```
In [44]: #With the reference of below link:
#https://www.kaggle.com/agungor2/various-confusion-matrix-plots
from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, optimal_model.predict(X_Test_Tfidf))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np
.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, f
mt='d')
```

Out[44]: <matplotlib.axes._subplots.AxesSubplot at 0x22683cfceb8>



[4.2.6]Classification Report

```
In [45]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction))
```

	precision	recall	f1-score	support
0	0.69	0.45	0.54	3010
1	0.59	0.80	0.68	2990
avg / total	0.64	0.62	0.61	6000

[4.3]Word2Vec

```
In [46]: i=0
list_of_sentence_train=[]
for sentence in X_Train:
    list_of_sentence_train.append(sentence.split())

w2v_model=Word2Vec(list_of_sentence_train,min_count=5,size=50, workers=4)
w2v_words = list(w2v_model.wv.vocab)
print("number of words that occurred minimum 5 times ",len(w2v_words))
print("sample words ", w2v_words[0:50])

number of words that occurred minimum 5 times 7526
sample words ['live', 'st', 'wanted', 'purchase', 'oh', 'henry', 'bars', 'co
uld', 'not', 'locate', 'market', 'bought', 'worried', 'might', 'melted', 'sum
mer', 'heat', 'no', 'issues', 'handed', 'son', 'born', 'everyone', 'loved',
'pleased', 'quality', 'pricing', 'item', 'use', 'homemade', 'rye', 'breads',
'would', 'become', 'unable', 'local', 'sources', 'caraway', 'fennel', 'seed
s', 'happy', 'discover', 'supplier', 'thought', 'try', 'flavor', 'k', 'cups',
'done', 'past']
```

[4.3.1] Compute avg w2v for train, test and CV data

```
In [47]: %%time
# average Word2Vec
# compute average word2vec for each review of train data.
sent_vectors_train = []; # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(list_of_sentence_train): # 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_train.append(sent_vec)
sent_vectors_train = np.array(sent_vectors_train)
print(sent_vectors_train.shape)
print(sent_vectors_train[0])
```

```
100%|███████████████████████████████████████████████████████████████████████████|
██████████ 14000/14000 [00:50<00:00, 277.54it/s]
```

```
(14000, 50)
[ 2.59178267e-02  2.56061523e-01 -1.79829784e-01 -5.69419847e-01
  4.26826637e-02  9.13998943e-02  3.90788424e-02  2.86464703e-01
  2.49460257e-01 -3.93666654e-01 -2.19688539e-01  3.56075176e-01
  1.24910339e-01 -2.37950270e-01 -9.33047577e-01  1.16164224e-01
  3.02268466e-01  3.93489492e-01  5.15138890e-01  1.24525340e-01
 -9.58608055e-02  3.14529083e-01 -1.59223592e-01  4.74623765e-01
 -6.41836816e-02 -3.89273322e-01  1.94994350e-01 -3.32118267e-01
  1.13790318e-01  9.24704226e-03 -3.25729935e-01  3.50075577e-04
  5.97130585e-01 -8.23682316e-01  6.28774566e-02  2.89672572e-01
  5.44256128e-02  3.55014000e-02 -5.68925405e-02  4.47809155e-01
 -1.31497630e-02  5.72693457e-03 -1.31890377e-02  4.19247218e-02
  2.38132485e-01  3.48509815e-01  4.72737611e-01  9.24839184e-02
  2.06765906e-01  1.89020746e-01]
```

Wall time: 50.8 s


```
In [48]: %%time
# compute average word2vec for each review of test data.
i=0
list_of_sentence_test=[]
for sentence in X_Test:
    list_of_sentence_test.append(sentence.split())

sent_vectors_test = []
for sent in tqdm(list_of_sentence_test):
    sent_vec = np.zeros(50)
    cnt_words = 0
    for word in sent:
        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%|██████████████████████████████████████████████████████████████████████████|  
██████████ 6000/6000 [00:48<00:00, 124.88it/s]
```

```
(6000, 50)
[ 0.17378092  0.30841917 -0.17879479 -0.31125192 -0.00700488  0.25553749
  0.28039384  0.50741675  0.2297549 -0.25229316 -0.23212694  0.23941967
  0.27791863  0.03789251 -1.00682962  0.11431529  0.3095299  0.14616974
  0.41085857 -0.02047079  0.0477207  0.49128097  0.02505071  0.47047255
  0.07435826 -0.71789804  0.17337335 -0.29543001 -0.04915047  0.06438296
 -0.20337037  0.16066263  0.65766505 -0.94801727 -0.13802876  0.75400569
 -0.08744782  0.30348756 -0.17039529  0.31141921  0.01799828  0.11887009
 -0.06853761  0.09746723  0.30995515  0.66880604  0.40245484  0.25822265
  0.10825151  0.02673325]
Wall time: 48.2 s
```

```
In [49]: %%time
# compute average word2vec for each review of cv data.
i=0
list_of_sentence_cv=[]
for sentence in X_cv:
    list_of_sentence_cv.append(sentence.split())

sent_vectors_cv = []
for sent in tqdm(list_of_sentence_cv):
    sent_vec = np.zeros(50)
    cnt_words = 0
    for word in sent:
        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%|███████████████████████████████████████████████████████████████████|
██████████6000/6000 [00:48<00:00, 123.14it/s]
```

```
(6000, 50)
[ 0.17378092  0.30841917 -0.17879479 -0.31125192 -0.00700488  0.25553749
  0.28039384  0.50741675  0.2297549  -0.25229316 -0.23212694  0.23941967
  0.27791863  0.03789251 -1.00682962  0.11431529  0.3095299  0.14616974
  0.41085857 -0.02047079  0.0477207  0.49128097  0.02505071  0.47047255
  0.07435826 -0.71789804  0.17337335 -0.29543001 -0.04915047  0.06438296
 -0.20337037  0.16066263  0.65766505 -0.94801727 -0.13802876  0.75400569
 -0.08744782  0.30348756 -0.17039529  0.31141921  0.01799828  0.11887009
 -0.06853761  0.09746723  0.30995515  0.66880604  0.40245484  0.25822265
  0.10825151  0.02673325]
Wall time: 48.8 s
```

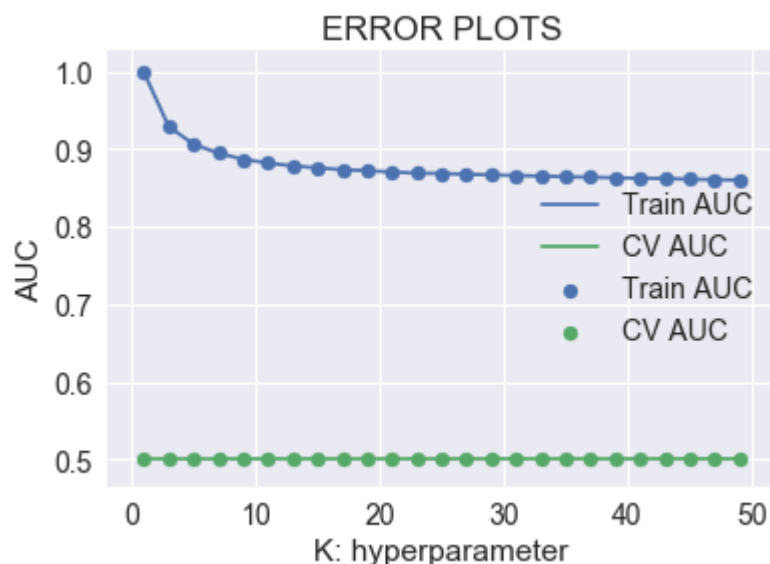
 [4.3.2] AUC Curve Plot

```

In [50]: train_AUC = []
CV_AUC = []
K = list(range(1,50,2))
for i in K:
    knn = KNeighborsClassifier(n_neighbors=i, algorithm='brute')
    #fit a model on Train TFIDF vectorizer
    knn.fit(sent_vectors_train, Y_Train)
    #predict probabilities on train TfIdf vectorizer
    Y_Train_Pred = knn.predict_proba(sent_vectors_train)[:,:1]
    #predict probabilities on cross validation TfIdf vectorizer
    Y_CV_Pred = knn.predict_proba(sent_vectors_cv)[:,:1]
    #calculate AUC score
    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()

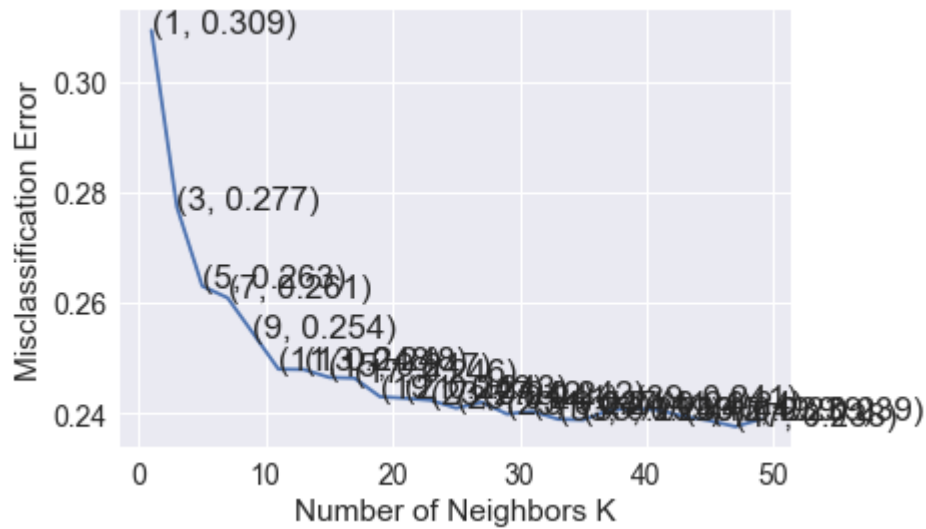
```



[4.3.3] 10-fold cross validation

```
In [51]: optimal_k3 = Optimal_K(sent_vectors_train, Y_Train)
print("optimal_k3:", optimal_k3)
```

The optimal number of neighbors is 47.



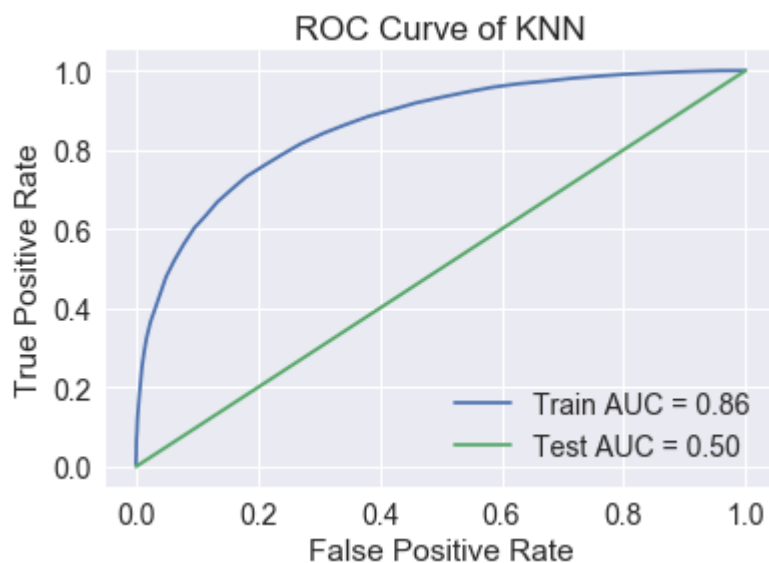
the misclassification error for each k value is : [0.309 0.277 0.263 0.261
0.254 0.248 0.248 0.247 0.246 0.243 0.243 0.242
0.241 0.242 0.24 0.24 0.239 0.239 0.24 0.241 0.24 0.239 0.239 0.238
0.239]
optimal_k3: 47

```
In [52]: optimal_model = KNeighborsClassifier(n_neighbors=optimal_k3, algorithm='brute'
)
optimal_model.fit(sent_vectors_train, Y_Train)
prediction = optimal_model.predict(sent_vectors_test)
```

[4.3.4]ROC Curve of KNN

```
In [53]: #with the reference of below link:
#https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn
#machine-learning-algorithm-using-python-and-sci
#predict probabilities on X_Train_Bow and X_Test_Bow and pass as param to roc_
#curve to find roc curve
Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, optimal_model.predict_proba(sent_vectors_train)[:,-1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, optimal_model.predict_proba(sent_vectors_test)[:,-1])
roc_auc4 = auc(Train_FPR, Train_TPR)
roc_auc5 = auc(Test_FPR, Test_TPR)

plt.title('Receiving Operating Characteristic Curve')
plt.plot(Train_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc4)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc5)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of KNN')
plt.show()
```



**[4.3.5]Train and Test Accuracy **

```
In [54]: Training_Accuracy = optimal_model.score(sent_vectors_train, Y_Train)
Training_Error = 1 - Training_Accuracy
Test_Accuracy = accuracy_score(Y_Test, prediction)
Test_Error = 1 - Test_Accuracy
```

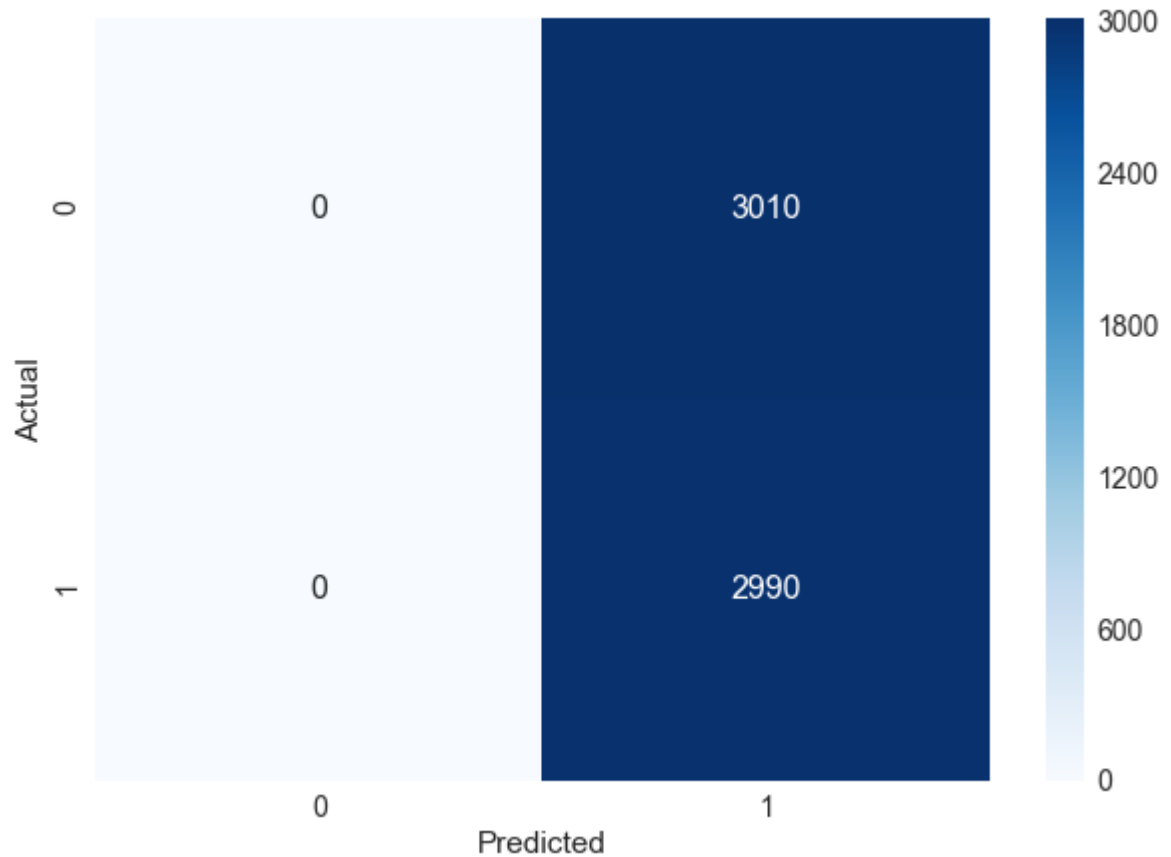
```
print("Training Accuracy: ", Training_Accuracy)
print("Train Error: ", Training_Error)
print("Test Accuracy: ", Test_Accuracy)
print("Test Error: ", Test_Error)
```

```
Training Accuracy:  0.7755714285714286
Train Error:  0.22442857142857142
Test Accuracy:  0.49833333333333335
Test Error:  0.5016666666666667
```

**[4.3.6]Confusion Matrix **

```
In [55]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, optimal_model.predict(sent_vectors_test))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, fmt='d')
```

Out[55]: <matplotlib.axes._subplots.AxesSubplot at 0x2268bcb1080>



[4.3.7]Classification Report

```
In [56]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction))
```

D:\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1135: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples.

```
'precision', 'predicted', average, warn_for)
```

	precision	recall	f1-score	support
0	0.00	0.00	0.00	3010
1	0.50	1.00	0.67	2990
avg / total	0.25	0.50	0.33	6000

[4.4]TF-IDF Weighted Word2Vec

```
In [57]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer()
model.fit(X_Train)
# 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_)))
```

[4.4.1] Compute the TF-IDF Weighted W2V for train, test and CV


```

In [58]: i=0
list_of_santance_train=[]
for sentence in X_Train:
    list_of_santance_train.append(sentence.split())
# 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_train = []; # the tfidf-w2v for each sentence/review is sto
red in this list
row=0;
for sent in tqdm(list_of_santance_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/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_model.wv 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

```

```

100%|████████████████████████████████████████████████████████████████████████████████|
████████████████████████████████████████████████████████████████████████████████ 14000/14000 [10:30<00:00, 4.25it/s]

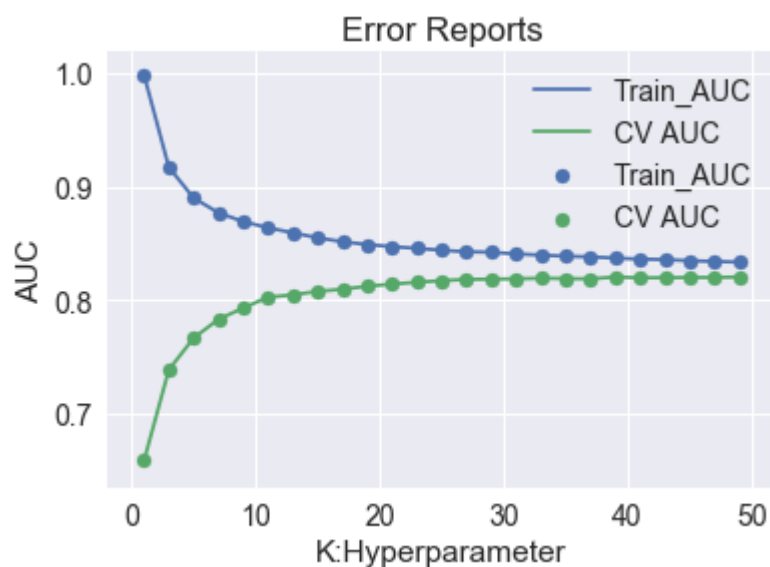
```



```

In [61]: train_AUC=[]
CV_AUC=[]
K = list(range(1,50,2))
for i in K:
    knn = KNeighborsClassifier(n_neighbors=i, algorithm='brute')
    #fit a model on Train TFIDF vectorizer
    knn.fit(tfidf_sent_vectors_train, Y_Train)
    #predict probabilities on train TfIdf vectorizer
    Y_Train_Pred = knn.predict_proba(tfidf_sent_vectors_train)[: ,1]
    #predict probabilities on cross validation TfIdf vectorizer
    Y_CV_Pred = knn.predict_proba(tfidf_sent_vectors_cv)[: ,1]
    #calculate AUC score
    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 Reports')
plt.show()

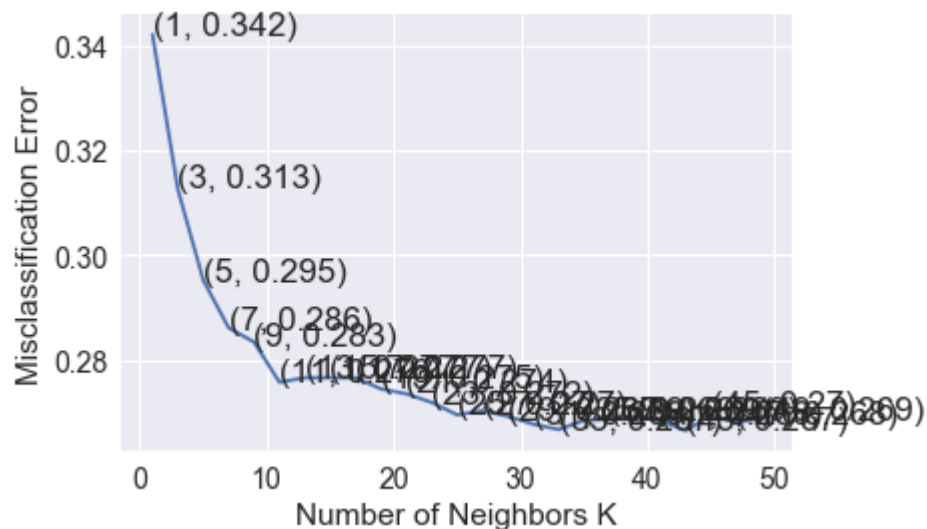
```



[4.4.3]10-fold cross validation

```
In [62]: optimal_k4 = Optimal_K(tfidf_sent_vectors_train, Y_Train)
print("optimal_k4:", optimal_k4)
```

The optimal number of neighbors is 43.



the misclassification error for each k value is : [0.342 0.313 0.295 0.286
0.283 0.276 0.277 0.277 0.277 0.275 0.274 0.272
0.27 0.27 0.269 0.268 0.267 0.269 0.269 0.269 0.268 0.267 0.27 0.268
0.269]
optimal_k4: 43

```
In [63]: optimal_model = KNeighborsClassifier(n_neighbors=optimal_k4, algorithm='brute'
)
optimal_model.fit(tfidf_sent_vectors_train, Y_Train)
prediction = optimal_model.predict(tfidf_sent_vectors_test)
```

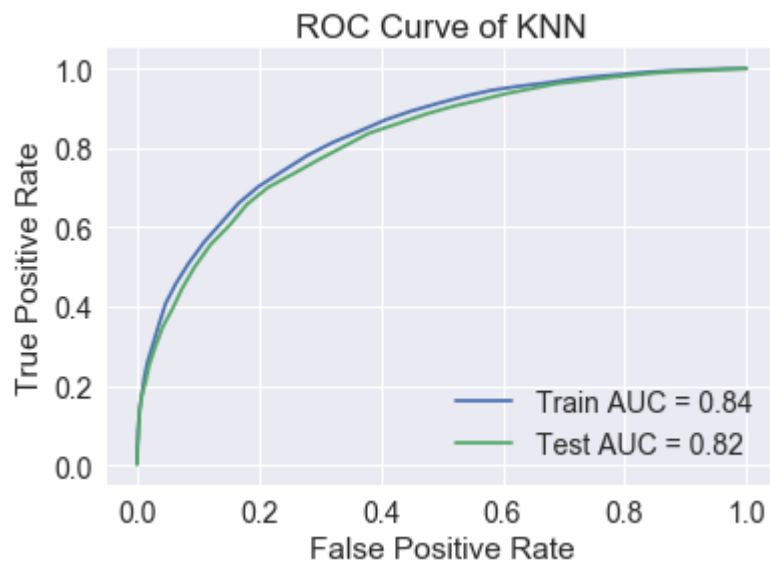
[4.4.4]ROC Curve of KNN

```

In [64]: #with the reference of below link:
#https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn
#machine-learning-algorithm-using-python-and-sci
#predict probabilities on X_Train_Bow and X_Test_Bow and pass as param to roc_
curve to find roc curve
Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, optimal_model.predict_proba(tfidf_sent_vectors_train)[:,-1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, optimal_model.predict_proba(tfidf_sent_vectors_test)[:,-1])
roc_auc6 = auc(Train_FPR, Train_TPR)
roc_auc7 = auc(Test_FPR, Test_TPR)

plt.title('Receiving Operating Characteristic Curve')
plt.plot(Train_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc6)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc7)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of KNN')
plt.show()

```



 [4.4.5]Train and Test Accuracy

```
In [65]: Training_Accuracy = optimal_model.score(tfidf_sent_vectors_train, Y_Train)
Training_Error = 1 - Training_Accuracy
Test_Accuracy = accuracy_score(Y_Test, prediction)
Test_Error = 1 - Test_Accuracy

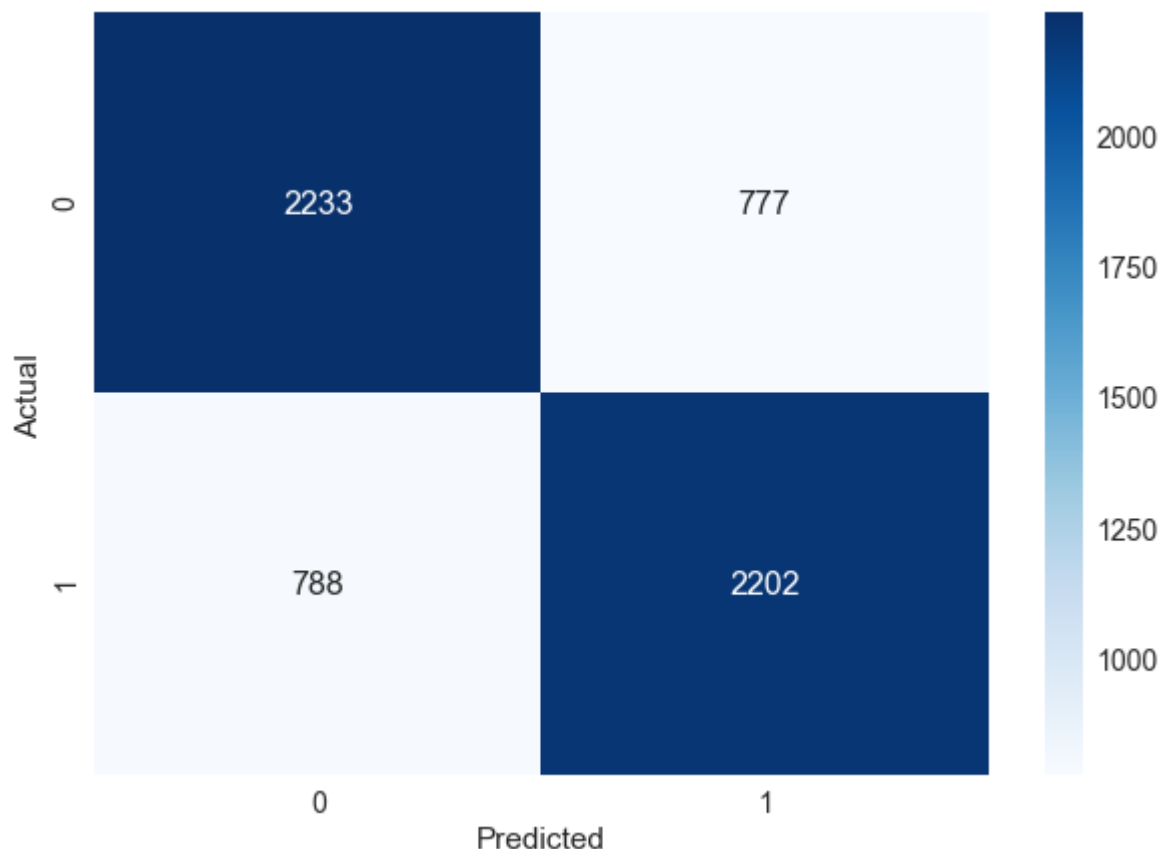
print("Training Accuracy: ", Training_Accuracy)
print("Train Error: ", Training_Error)
print("Test Accuracy: ", Test_Accuracy)
print("Test Error: ", Test_Error)
```

```
Training Accuracy:  0.7514285714285714
Train Error:  0.24857142857142855
Test Accuracy:  0.7391666666666666
Test Error:  0.26083333333333336
```

[4.4.6]Confusion Matrix

```
In [66]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, optimal_model.predict(tfidf_sent_vectors_test))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, fmt='d')
```

Out[66]: <matplotlib.axes._subplots.AxesSubplot at 0x2268bc9e470>



[4.4.7]Classification Report

```
In [67]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction))
```

	precision	recall	f1-score	support
0	0.74	0.74	0.74	3010
1	0.74	0.74	0.74	2990
avg / total	0.74	0.74	0.74	6000


```
In [68]: data_pos1 = final[final["Score"] == 1].sample(n = 1000)
data_neg1 = final[final["Score"] == 0].sample(n = 1000)
final2 = pd.concat([data_pos1, data_neg1])
final2.shape
```

```
Out[68]: (2000, 11)
```

```
In [69]: Y = final2['Score'].values
X = final2['cleaned_text'].values
print(Y.shape)
print(type(Y))
print(X.shape)
print(type(X))
```

```
(2000,)
<class 'numpy.ndarray'>
(2000,)
<class 'numpy.ndarray'>
```

 Split the data set into train, test and CV

```
In [70]: # split the data set into train and test
X_Train, X_Test, Y_Train, Y_Test = train_test_split(X,Y,test_size=0.2, random_
state=12)

# split the train data set into cross validation train and cross validation te
st
X_tr, X_cv, Y_tr, Y_cv = train_test_split(X,Y, test_size=0.2, random_state=12)

#X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.2,random_state
=12)
#X_train,X_cv,Y_train,Y_cv=train_test_split(X,Y,test_size=0.2,random_state=12)
print('='*100)
print("After splitting")
print("X_Train Shape:",X_Train.shape,"Y_Train Shape:",Y_Train.shape)
print("X_cv Shape:",X_cv.shape,"Y_cv Shape",Y_cv.shape)
print("X_Test Shape",X_Test.shape,"Y_Test Shape",Y_Test.shape)

=====
=====
After splitting
X_Train Shape: (1600,) Y_Train Shape: (1600,)
X_cv Shape: (400,) Y_cv Shape (400,)
X_Test Shape (400,) Y_Test Shape (400,)
```

Featurization - Using KD-Tree Algorithm

[5.1]Bag of Words

```
In [71]: #Bow
count_vect = CountVectorizer(min_df=10, max_features=500) #in scikit-learn
count_vect.fit(X_Train)
print("some feature names ", count_vect.get_feature_names()[:10])
X_Train_Bow_kdtree = count_vect.transform(X_Train)
X_Test_Bow_kdtree = count_vect.transform(X_Test)
X_CV_Bow_kdtree = count_vect.transform(X_cv)

print('='*50)

#final_counts = count_vect.transform(X_Test)

print("the type of X Train : ",type(X_Train_Bow_kdtree))
print("the shape of Train BOW vectorizer ",X_Train_Bow_kdtree.get_shape())
print("the shape of Test BOW vectorizer ",X_Test_Bow_kdtree.get_shape())
print("the shape of CV BOW vectorizer ",X_CV_Bow_kdtree.get_shape())
#print("the number of unique words ", final_counts.get_shape()[1])

some feature names  ['able', 'absolutely', 'actually', 'add', 'added', 'addin
g', 'aftertaste', 'ago', 'almost', 'along']
=====
the type of X Train :  <class 'scipy.sparse.csr.csr_matrix'>
the shape of Train BOW vectorizer  (1600, 500)
the shape of Test BOW vectorizer  (400, 500)
the shape of CV BOW vectorizer  (400, 500)
```

[5.1.1]AUC Curve Plot

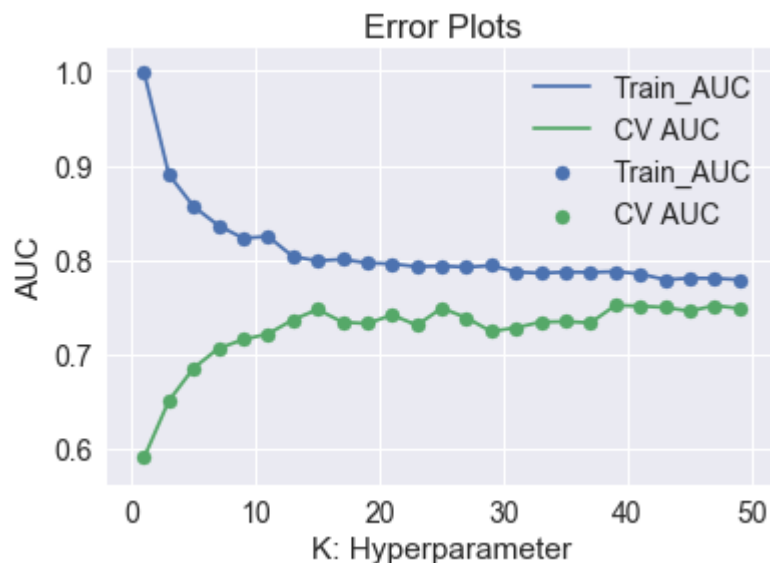
```

In [72]: train_AUC = []
        CV_AUC = []

        K = list(range(1,50,2))
        for i in K:
            knn = KNeighborsClassifier(n_neighbors = i, algorithm = 'kd_tree')
            #fit a model on train BOW vectorizer
            knn.fit(X_Train_Bow_kdtree.todense(), Y_Train)
            #predict probabilities on train BOW vectorizer
            Y_Train_Pred = knn.predict_proba(X_Train_Bow_kdtree.todense())[:,1]
            #predict probabilities on Cross validation BOW vectorizer
            Y_CV_Pred = knn.predict_proba(X_CV_Bow_kdtree.todense())[:,1]
            #calculate AUC score
            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()

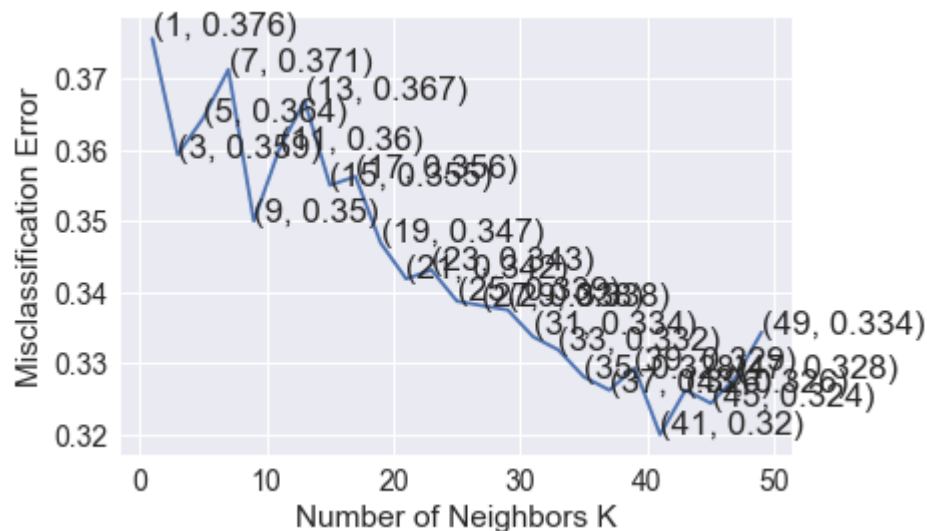
```



**[5.1.2]10-fold cross validation **

```
In [73]: optimal_k1_kdtree = Optimal_K(X_Train_Bow_kdtree.todense(), Y_Train)
print("optimal_k1_kdtree:", optimal_k1_kdtree)
```

The optimal number of neighbors is 41.



```
In [74]: optimal_model = KNeighborsClassifier(n_neighbors=optimal_k1_kdtree, algorithm=
'kd_tree')
optimal_model.fit(X_Train_Bow_kdtree, Y_Train)
prediction = optimal_model.predict(X_Test_Bow_kdtree)
```

the misclassification error for each k value is : [0.376 0.359 0.364 0.371
0.35 0.36 0.367 0.355 0.356 0.347 0.342 0.343
0.339 0.338 0.338 0.334 0.332 0.328 0.326 0.329 0.32 0.326 0.324 0.328
0.334]
optimal_k1_kdtree: 41

D:\Anaconda3\lib\site-packages\sklearn\neighbors\base.py:212: UserWarning: cannot use tree with sparse input: using brute force
warnings.warn("cannot use tree with sparse input: ")

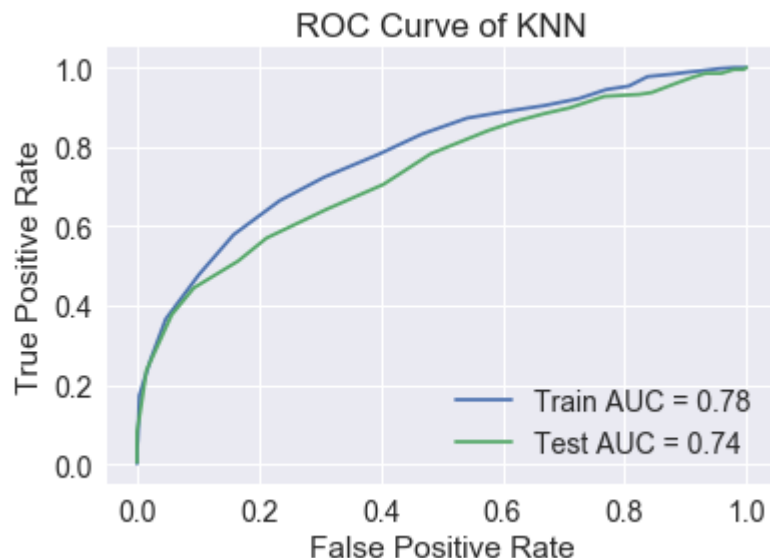
[5.1.3]ROC Curve of KNN

```

In [75]: #with the reference of below link:
#https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn
#machine-learning-algorithm-using-python-and-sci
#predict probabilities on X_Train_Bow and X_Test_Bow and pass as param to roc_
#curve to find roc curve
Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, optimal_model.predict_proba(X_Train_Bow_kdtree)[:,-1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, optimal_model.predict_proba(X_Test_Bow_kdtree)[:,-1])
roc_auc8 = auc(Train_FPR, Train_TPR)
roc_auc9 = auc(Test_FPR, Test_TPR)

plt.title('Receiving Operating Characteristic Curve')
plt.plot(Train_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc8)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc9)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of KNN')
plt.show()

```



[5.1.4]Train and Test Accuracy

```
In [76]: Training_Accuracy = optimal_model.score(X_Train_Bow_kdtree, Y_Train)
print('Training_Accuracy=%0.3f'%Training_Accuracy)
Training_Error = 1 - Training_Accuracy
print('Training_Error=%0.3f'%Training_Error)

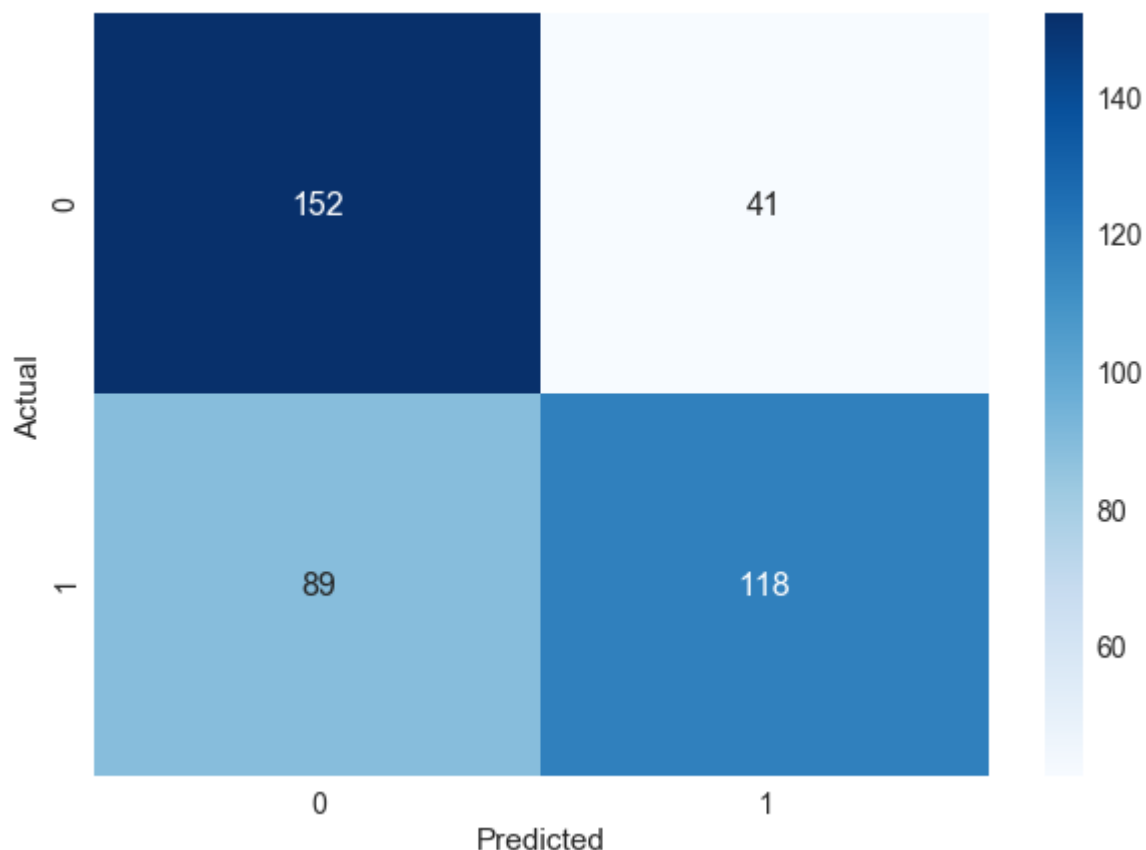
Test_Accuracy = accuracy_score(Y_Test, prediction)
print('Test_Accuracy=%0.3f'%Test_Accuracy)
Test_Error = 1 - Test_Accuracy
print('Test_Error=%0.3f'%Test_Error)

Training_Accuracy=0.716
Training_Error=0.284
Test_Accuracy=0.675
Test_Error=0.325
```

[5.1.5]Confusion Matrix

```
In [77]: #With the reference of below link:
#https://www.kaggle.com/agungor2/various-confusion-matrix-plots
from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, optimal_model.predict(X_Test_Bow_kdtree
))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np
.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, f
mt='d')
```

Out[77]: <matplotlib.axes._subplots.AxesSubplot at 0x2268bc985c0>



[5.1.6]Classification Report

```
In [78]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction))
```

	precision	recall	f1-score	support
0	0.63	0.79	0.70	193
1	0.74	0.57	0.64	207
avg / total	0.69	0.68	0.67	400

[5.2]TF-IDF

```
In [79]: #TF-IDF
tf_idf_vect_kdtree = TfidfVectorizer(ngram_range=(1,2), min_df=5, max_df=100)
tf_idf_vect_kdtree.fit_transform(X_Train)
print("some sample features(unique words in the corpus)",tf_idf_vect_kdtree.get_feature_names()[0:10])
print('='*50)

X_Train_TfIdf_kdtree = count_vect.transform(X_Train)
X_Test_TfIdf_kdtree = count_vect.transform(X_Test)
X_CV_TfIdf_kdtree = count_vect.transform(X_cv)

#final_tf_idf = tf_idf_vect.transform(X_Test)
print("the type of count vectorizer ",type(X_Train_TfIdf_kdtree))
print("the shape of out text TFIDF vectorizer ",X_Train_TfIdf_kdtree.get_shape())
print("the shape of out text TFIDF vectorizer ",X_Test_TfIdf_kdtree.get_shape())
print("the shape of out text TFIDF vectorizer ",X_CV_TfIdf_kdtree.get_shape())
#print("the number of unique words including both unigrams and bigrams ", final_tf_idf.get_shape()[1])

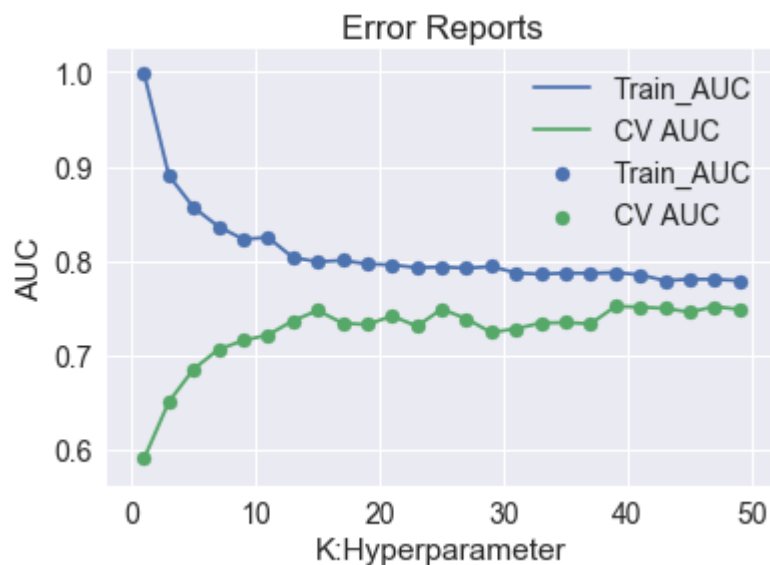
some sample features(unique words in the corpus) ['ability', 'able', 'able find', 'absolute', 'absolutely', 'absolutely delicious', 'absolutely love', 'accept', 'according', 'acid']
=====
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer (1600, 500)
the shape of out text TFIDF vectorizer (400, 500)
the shape of out text TFIDF vectorizer (400, 500)
```

[5.2.1]AUC Curve Plot


```

In [80]: train_AUC=[]
CV_AUC=[]
K = list(range(1,50,2))
for i in K:
    knn = KNeighborsClassifier(n_neighbors=i, algorithm='kd_tree')
    #fit a model on Train TFIDF vectorizer
    knn.fit(X_Train_TfIdf_kdtree.todense(), Y_Train)
    #predict probabilities on train TfIdf vectorizer
    Y_Train_Pred = knn.predict_proba(X_Train_TfIdf_kdtree.todense())[:,1]
    #predict probabilities on cross validation TfIdf vectorizer
    Y_CV_Pred = knn.predict_proba(X_CV_TfIdf_kdtree.todense())[:,1]
    #calculate AUC score
    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 Reports')
plt.show()

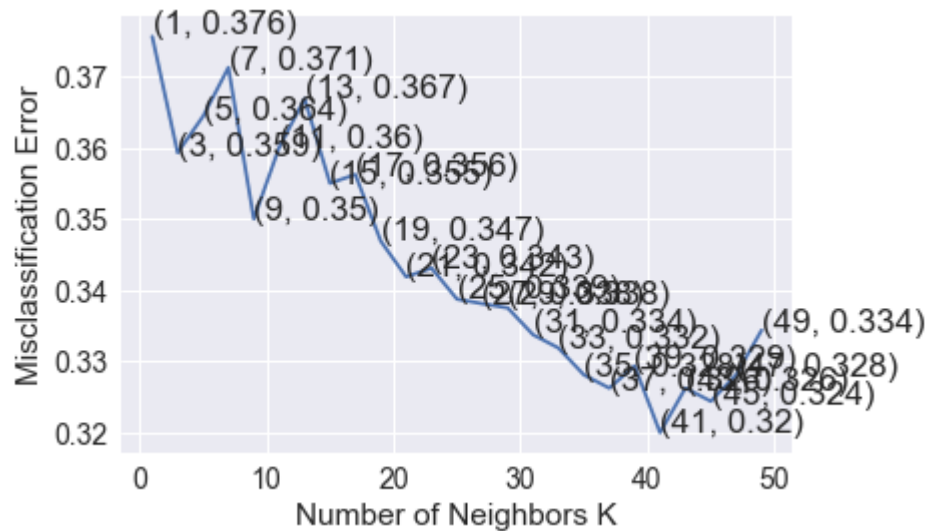
```



**[5.2.2]10-fold cross validation **

```
In [81]: optimal_k2_kdtree = Optimal_K(X_Train_TfIdf_kdtree.todense(), Y_Train)
print("optimal_k2_kdtree:", optimal_k2_kdtree)
```

The optimal number of neighbors is 41.



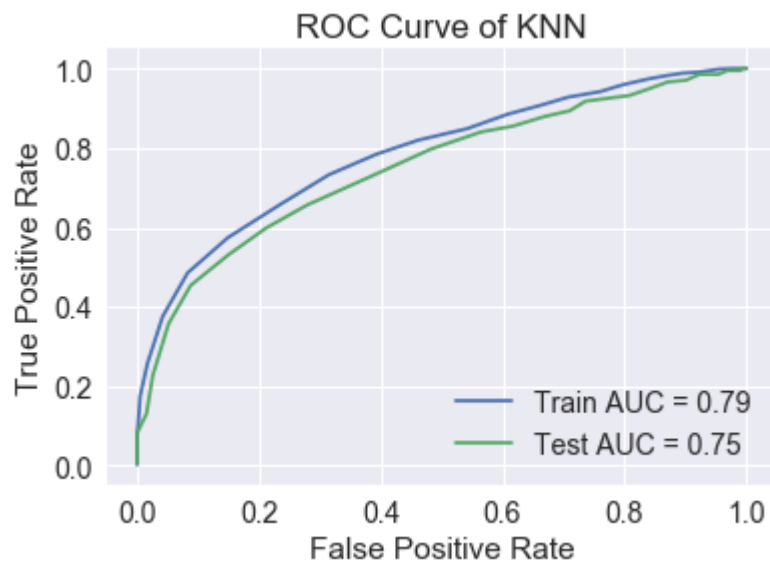
```
the misclassification error for each k value is : [0.376 0.359 0.364 0.371
0.35 0.36 0.367 0.355 0.356 0.347 0.342 0.343
0.339 0.338 0.338 0.334 0.332 0.328 0.326 0.329 0.32 0.326 0.324 0.328
0.334]
optimal_k2_kdtree: 41
```

```
In [82]: optimal_model = KNeighborsClassifier(n_neighbors=optimal_k2_kdtree, algorithm=
'kd_tree')
optimal_model.fit(X_Train_TfIdf_kdtree.todense(), Y_Train)
prediction = optimal_model.predict(X_Test_TfIdf_kdtree.todense())
```

[5.2.3]ROC Curve of KNN

```
In [83]: #with the reference of below link:
#https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn
#machine-learning-algorithm-using-python-and-sci
#predict probabilities on X_Train_Bow and X_Test_Bow and pass as param to roc_
curve to find roc curve
Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, optimal_model.predict_proba(X_Train_Tfidf_kdtree.todense()))[:,1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, optimal_model.predict_proba(X_Test_Tfidf_kdtree.todense()))[:,1])
roc_auc10 = auc(Ttrain_FPR, Train_TPR)
roc_auc11 = auc(Test_FPR, Test_TPR)

plt.title('Receiving Operating Characteristic Curve')
plt.plot(Ttrain_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc10)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc11)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of KNN')
plt.show()
```



[5.2.4]Train and Test Accuracy

```
In [84]: Training_Accuracy = optimal_model.score(X_Train_TfIdf_kdtree.todense(), Y_Train)
print('Training_Accuracy=%0.3f'%Training_Accuracy)
Training_Error = 1 - Training_Accuracy
print('Training_Error=%0.3f'%Training_Error)

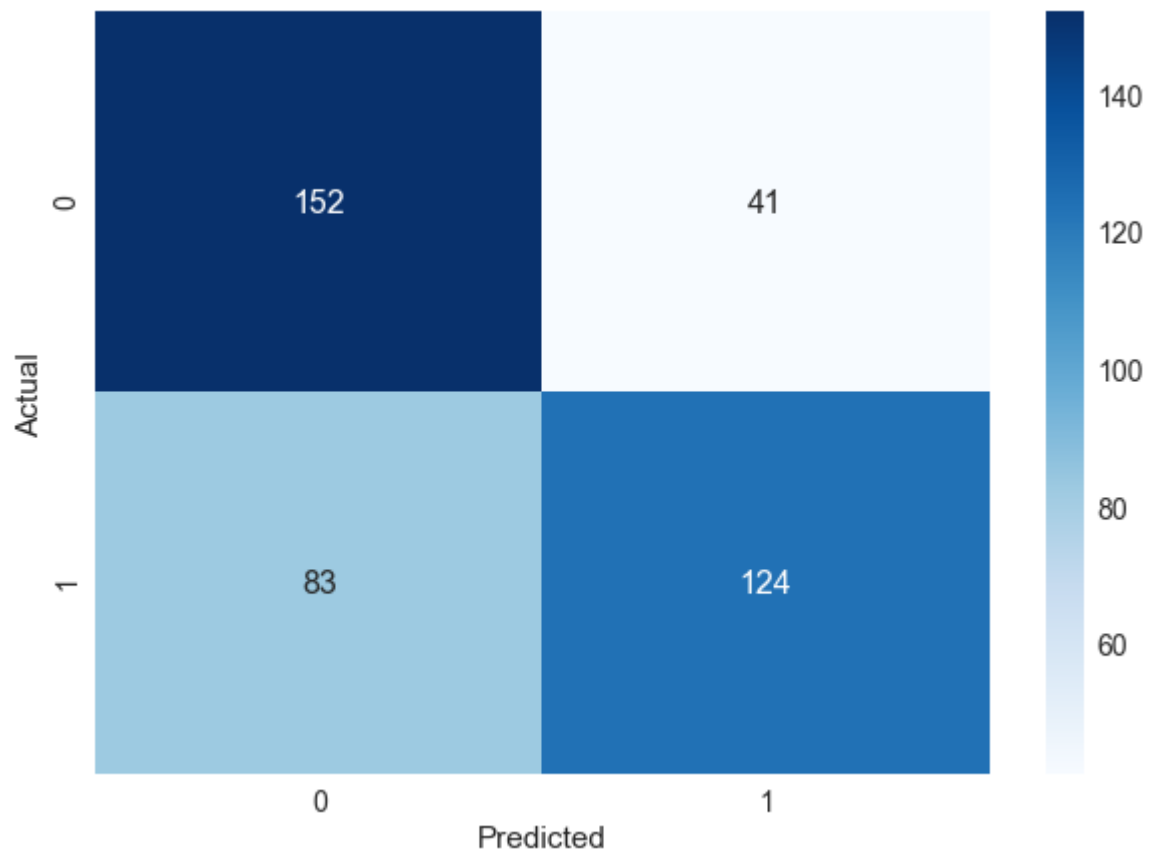
Test_Accuracy = accuracy_score(Y_Test, prediction)
print('Test_Accuracy=%0.3f'%Test_Accuracy)
Test_Error = 1 - Test_Accuracy
print('Test_Error=%0.3f'%Test_Error)

Training_Accuracy=0.712
Training_Error=0.288
Test_Accuracy=0.690
Test_Error=0.310
```

[5.2.5]Confusion Matrix

```
In [85]: #With the reference of below link:
#https://www.kaggle.com/agungor2/various-confusion-matrix-plots
from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, optimal_model.predict(X_Test_Tfidf_kdtr
ee.todense()))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np
.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, f
mt='d')
```

Out[85]: <matplotlib.axes._subplots.AxesSubplot at 0x2268bc98e10>



[5.2.6]Classification Report

```
In [86]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction))
```

	precision	recall	f1-score	support
0	0.65	0.79	0.71	193
1	0.75	0.60	0.67	207
avg / total	0.70	0.69	0.69	400

[5.3]Word2Vec

```
In [87]: i=0
list_of_senatnce_train=[]
for sentence in X_Train:
    list_of_senatnce_train.append(sentence.split())

w2v_model=Word2Vec(list_of_senatnce_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])

number of words that occured minimum 5 times 2169
sample words ['bought', 'try', 'different', 'flavor', 'filled', 'hazelnut',
'instead', 'good', 'ordered', 'arrived', 'timely', 'didnt', 'want', 'go', 're
turning', 'despite', 'expiration', 'date', 'nuts', 'rancid', 'tasty', 'produc
t', 'not', 'k', 'cup', 'coffee', 'tasted', 'like', 'water', 'two', 'year', 'o
ld', 'trying', 'everything', 'keep', 'teeth', 'clean', 'avoid', 'putting', 'g
uys', 'would', 'near', 'look', 'great', 'smell', 'finally', 'gave', 'rescue',
'pup', 'still']
```

[5.3.1] Computer avg w2v for train, test, and CV

```
100%|██████████████████████████████████████████████████████████████████████████|  
██████████ 1600/1600 [00:04<00:00, 383.15it/s]
```

```
(1600, 50)
[ 3.92335101e-01  2.54011220e-01 -8.90799452e-02 -5.46659711e-01
  7.44760209e-02 -2.30371141e-02  5.21362415e-01  3.14892812e-01
  2.36234717e-01 -2.57891757e-01 -1.18026758e-01  5.20183303e-01
  1.93848384e-02 -3.68769039e-01 -7.03955669e-01 -2.71730401e-01
  7.17273440e-01  1.92535790e-01  6.46719663e-01  8.02793867e-02
 -1.80937334e-02  1.05153423e-01 -2.38693434e-02  3.17114180e-01
  1.63557555e-01 -7.80010068e-01  4.53045098e-01 -5.06500697e-01
 -8.38716228e-02 -1.99799498e-02 -2.02212512e-01 -2.02313494e-01
  3.97812624e-01 -1.13988541e+00  5.71482819e-02  6.95349629e-01
 -6.19616081e-02  1.23181971e-01 -1.02574336e-01  2.99539909e-01
 -5.78120245e-01 -1.94581533e-01 -6.66226259e-02  5.85843971e-04
  3.57824917e-01  3.71588803e-01  6.22277365e-01  2.92814677e-01
  1.12035273e-02 -5.98712191e-02]
```

```
In [89]: %%time

i=0
list_of_sentence_test=[]
for sentence in X_Test:
    list_of_sentence_test.append(sentence.split())

# 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_sentence_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])
```

```
100%|██████████████████████████████████████████████████████████████████████████|
██████████ | 400/400 [00:00<00:00, 490.26it/s]
```

```
(400, 50)
[ 0.41254308  0.2666285 -0.09195531 -0.57477557  0.078032 -0.02571928
  0.54536079  0.32623301  0.25118605 -0.27070716 -0.12019518  0.54836249
  0.02067853 -0.3870541 -0.73622064 -0.28655878  0.75471893  0.20268793
  0.6785012  0.08363406 -0.01952433  0.1105312 -0.02893279  0.32983514
  0.17297649 -0.81662611  0.47489748 -0.53308903 -0.09103425 -0.01711964
 -0.21310112 -0.21276142  0.42268098 -1.19668647  0.06092215  0.7278375
 -0.0652902  0.12765201 -0.10990935  0.31306563 -0.60916663 -0.20499381
 -0.0703765  0.00247578  0.3752465  0.38964954  0.64951533  0.30290898
  0.01304127 -0.06000011]
```



```
In [90]: %%time
i=0
list_of_sentence_cv=[]
for sentence in X_cv:
    list_of_sentence_cv.append(sentence.split())
# 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_sentence_cv): # 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_cv.append(sent_vec)
sent_vectors_cv = np.array(sent_vectors_cv)
print(sent_vectors_cv.shape)
print(sent_vectors_cv[0])
```

```
100%|██████████████████████████████████████████████████████████████████████████|
██████████ | 400/400 [00:00<00:00, 511.99it/s]
```

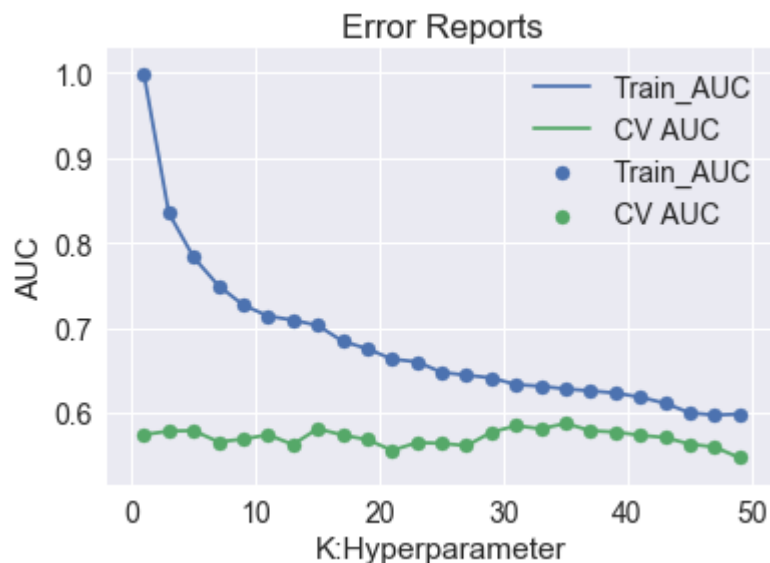
```
(400, 50)
[ 0.41254308  0.2666285 -0.09195531 -0.57477557  0.078032 -0.02571928
  0.54536079  0.32623301  0.25118605 -0.27070716 -0.12019518  0.54836249
  0.02067853 -0.3870541 -0.73622064 -0.28655878  0.75471893  0.20268793
  0.6785012  0.08363406 -0.01952433  0.1105312 -0.02893279  0.32983514
  0.17297649 -0.81662611  0.47489748 -0.53308903 -0.09103425 -0.01711964
 -0.21310112 -0.21276142  0.42268098 -1.19668647  0.06092215  0.7278375
 -0.0652902  0.12765201 -0.10990935  0.31306563 -0.60916663 -0.20499381
 -0.0703765  0.00247578  0.3752465  0.38964954  0.64951533  0.30290898
  0.01304127 -0.06000011]
```

[5.3.2] AUC plot

```

In [92]: train_AUC=[]
CV_AUC=[]
K = list(range(1,50,2))
for i in K:
    knn = KNeighborsClassifier(n_neighbors=i, algorithm='kd_tree')
    #fit a model on Train TFIDF vectorizer
    knn.fit(sent_vectors_train, Y_Train)
    #predict probabilities on train TfIdf vectorizer
    Y_Train_Pred = knn.predict_proba(sent_vectors_train)[: ,1]
    #predict probabilities on cross validation TfIdf vectorizer
    Y_CV_Pred = knn.predict_proba(sent_vectors_cv)[: ,1]
    #calculate AUC score
    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 Reports')
plt.show()

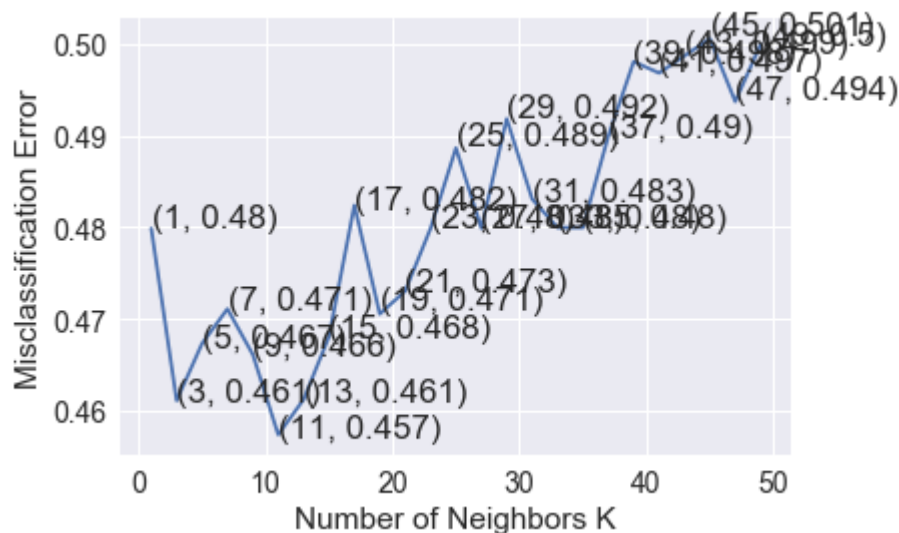
```



**[5.3.3]10-fold cross validation **

```
In [94]: optimal_k3_kdtree = Optimal_K(sent_vectors_train, Y_Train)
print("optimal_k3_kdtree:", optimal_k3_kdtree)
```

The optimal number of neighbors is 11.



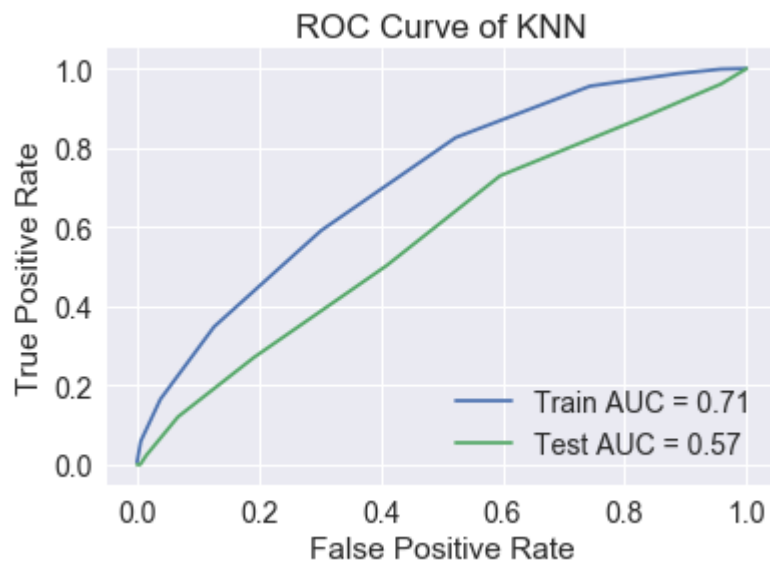
the misclassification error for each k value is : [0.48 0.461 0.467 0.471
0.466 0.457 0.461 0.468 0.482 0.471 0.473 0.48
0.489 0.48 0.492 0.483 0.48 0.48 0.49 0.498 0.497 0.499 0.501 0.494
0.5]
optimal_k3_kdtree: 11

```
In [95]: optimal_model = KNeighborsClassifier(n_neighbors=optimal_k3_kdtree, algorithm=
'kd_tree')
optimal_model.fit(sent_vectors_train, Y_Train)
prediction = optimal_model.predict(sent_vectors_test)
```

[5.3.4]ROC Curve of KNN

```
In [96]: #with the reference of below link:
#https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn
#machine-learning-algorithm-using-python-and-sci
#predict probabilities on X_Train_Bow and X_Test_Bow and pass as param to roc_
#curve to find roc curve
Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, optimal_model.predict_proba(sent_vectors_train)[:,-1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, optimal_model.predict_proba(sent_vectors_test)[:,-1])
roc_auc12 = auc(Train_FPR, Train_TPR)
roc_auc13 = auc(Test_FPR, Test_TPR)

plt.title('Receiving Operating Characteristic Curve')
plt.plot(Train_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc12)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc13)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of KNN')
plt.show()
```



**[5.3.5]Train and Test Accuracy **

```
In [97]: Training_Accuracy = optimal_model.score(sent_vectors_train, Y_Train)
Training_Error = 1 - Training_Accuracy
Test_Accuracy = accuracy_score(Y_Test, prediction)
Test_Error = 1 - Test_Accuracy

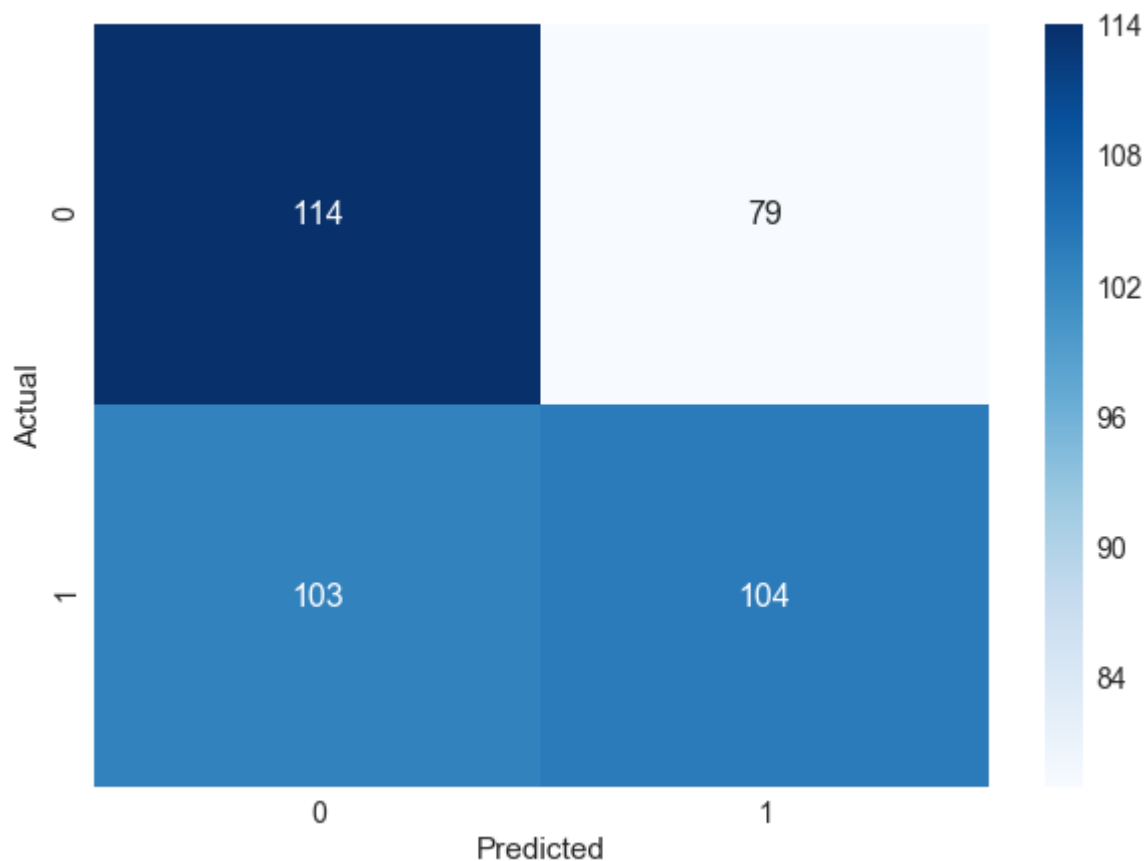
print("Training Accuracy: ", Training_Accuracy)
print("Train Error: ", Training_Error)
print("Test Accuracy: ", Test_Accuracy)
print("Test Error: ", Test_Error)
```

```
Training Accuracy:  0.645
Train Error:  0.355
Test Accuracy:  0.545
Test Error:  0.45499999999999996
```

**[5.3.6]Confusion Matrix **

```
In [98]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, optimal_model.predict(sent_vectors_test))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, fmt='d')
```

Out[98]: <matplotlib.axes._subplots.AxesSubplot at 0x226811a1470>



[5.3.7]Classification Report

```
In [99]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction))
```

	precision	recall	f1-score	support
0	0.53	0.59	0.56	193
1	0.57	0.50	0.53	207
avg / total	0.55	0.55	0.54	400

[5.4]TF-IDF

```
In [100]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
          model = TfidfVectorizer()
          model.fit(X_Train)
          # 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_)))
```

[5.4.1] Compute TF-IDF weighted Word2Vec for Train, Test, and CV

```
In [101]: i=0
          list_of_sentence_train=[]
          for sentence in X_Train:
              list_of_sentence_train.append(sentence.split())
          # 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_train = []; # the tfidf-w2v for each sentence/review is sto
            red in this list
          row=0;
          for sent in tqdm(list_of_sentence_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/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_sent_vectors_train.append(sent_vec)
              row += 1
```

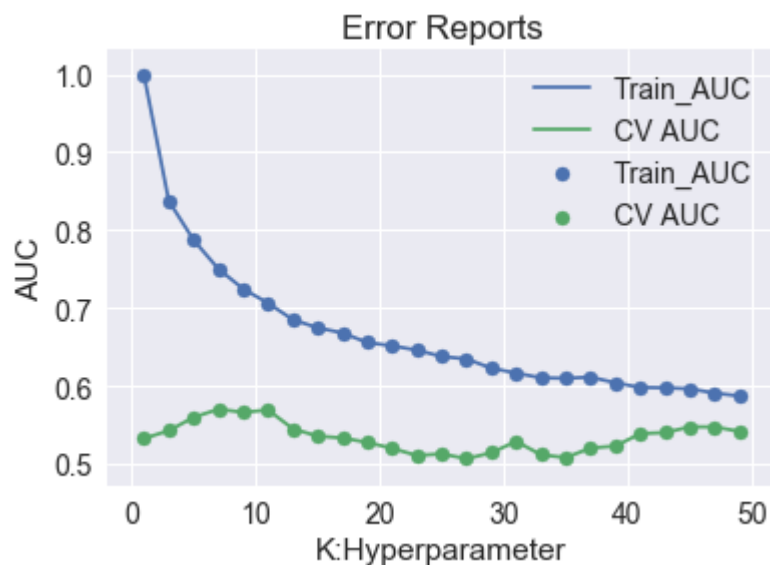
```
100%|████████████████████████████████████████████████████████████████████████████████|
████████████████████████████████████████████████████████████████████████████████| 1600/1600 [00:13<00:00, 81.00it/s]
```



```

In [105]: train_AUC=[]
CV_AUC=[]
K = list(range(1,50,2))
for i in K:
    knn = KNeighborsClassifier(n_neighbors=i, algorithm='kd_tree')
    #fit a model on Train TFIDF vectorizer
    knn.fit(tfidf_sent_vectors_train, Y_Train)
    #predict probabilities on train TfIdf vectorizer
    Y_Train_Pred = knn.predict_proba(tfidf_sent_vectors_train)[:,:1]
    #predict probabilities on cross validation TfIdf vectorizer
    Y_CV_Pred = knn.predict_proba(tfidf_sent_vectors_cv)[:,:1]
    #calculate AUC score
    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 Reports')
plt.show()

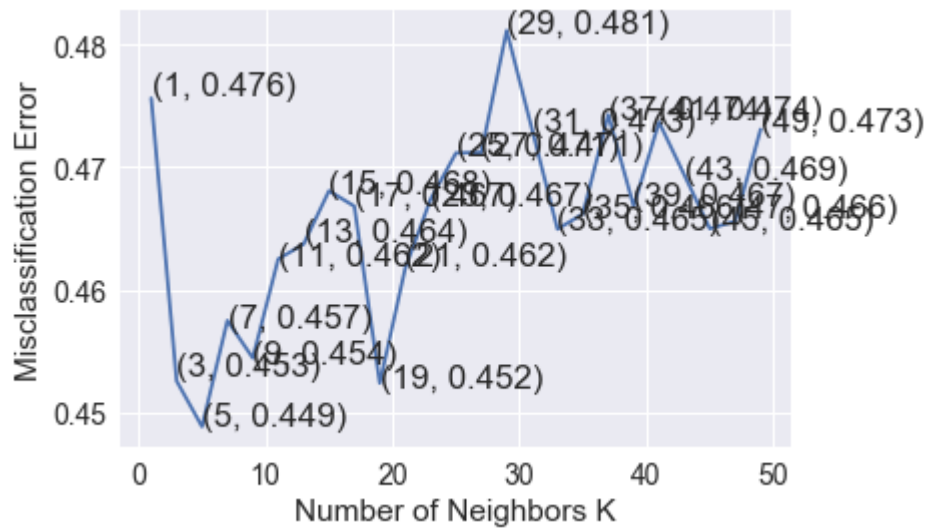
```



**[5.4.3]10-fold cross validation **

```
In [106]: optimal_k4_kdtree = Optimal_K(tfidf_sent_vectors_train, Y_Train)
print("optimal_k4_kdtree:", optimal_k4_kdtree)
```

The optimal number of neighbors is 5.



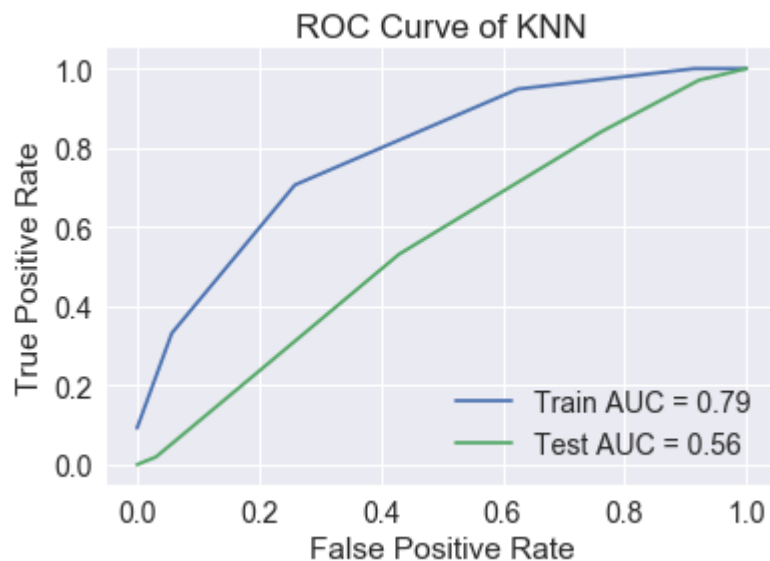
the misclassification error for each k value is : [0.476 0.453 0.449 0.457
0.454 0.462 0.464 0.468 0.467 0.452 0.462 0.467
0.471 0.471 0.481 0.473 0.465 0.466 0.474 0.467 0.474 0.469 0.465 0.466
0.473]
optimal_k4_kdtree: 5

```
In [107]: optimal_model = KNeighborsClassifier(n_neighbors=optimal_k4_kdtree, algorithm=
'kd_tree')
optimal_model.fit(tfidf_sent_vectors_train, Y_Train)
prediction = optimal_model.predict(tfidf_sent_vectors_test)
```

[5.4.4] ROC Curve of KNN

```
In [108]: #with the reference of below link:
#https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn
#machine-learning-algorithm-using-python-and-sci
#predict probabilities on X_Train_Bow and X_Test_Bow and pass as param to roc_
#curve to find roc curve
Train_FPR, Train_TPR, threshold = roc_curve(Y_Train, optimal_model.predict_proba(tfidf_sent_vectors_train)[:,-1])
Test_FPR, Test_TPR, threshold = roc_curve(Y_Test, optimal_model.predict_proba(tfidf_sent_vectors_test)[:,-1])
roc_auc14 = auc(Ttrain_FPR, Train_TPR)
roc_auc15 = auc(Test_FPR, Test_TPR)

plt.title('Receiving Operating Characteristic Curve')
plt.plot(Ttrain_FPR, Train_TPR, label = 'Train AUC = %0.2f' % roc_auc14)
plt.plot(Test_FPR, Test_TPR, label = 'Test AUC = %0.2f' % roc_auc15)
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of KNN')
plt.show()
```



[5.4.5]Train and Test Accuracy

```
In [109]: Training_Accuracy = optimal_model.score(tfidf_sent_vectors_train, Y_Train)
Training_Error = 1 - Training_Accuracy
Test_Accuracy = accuracy_score(Y_Test, prediction)
Test_Error = 1 - Test_Accuracy

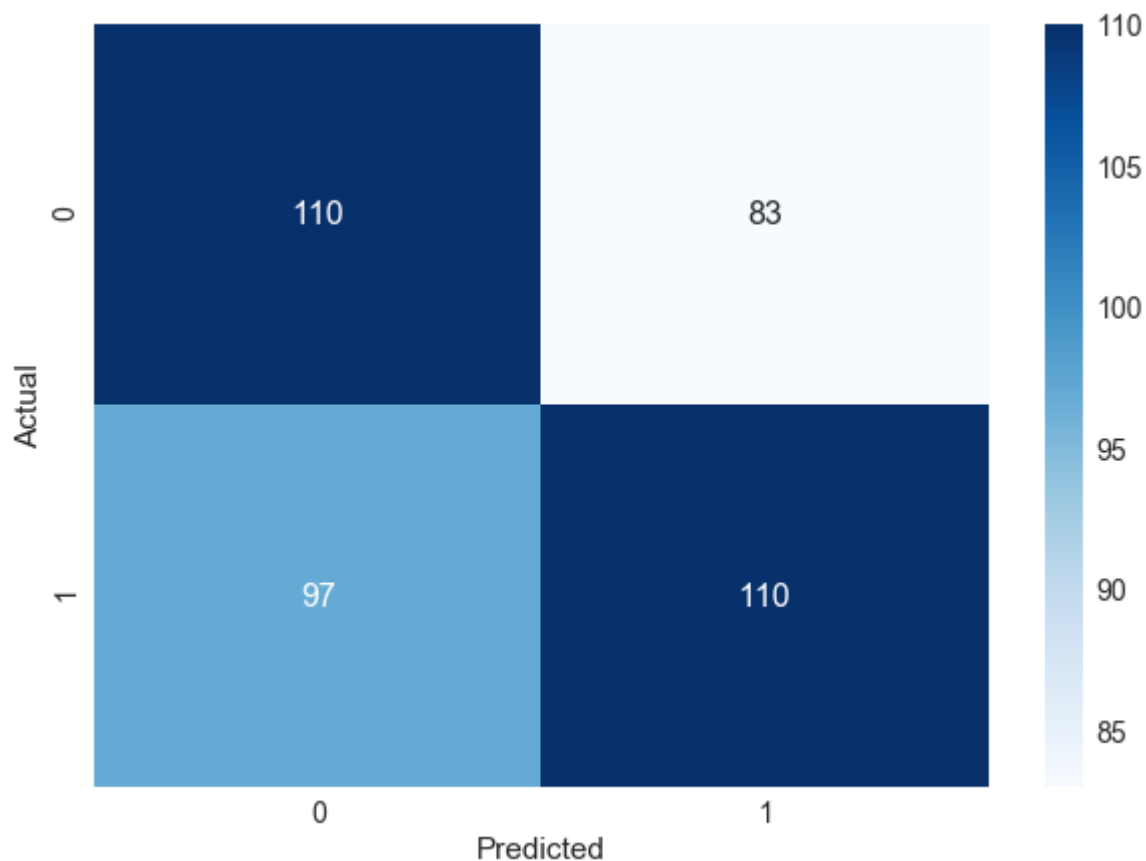
print("Training Accuracy: ", Training_Accuracy)
print("Train Error: ", Training_Error)
print("Test Accuracy: ", Test_Accuracy)
print("Test Error: ", Test_Error)
```

```
Training Accuracy:  0.72375
Train Error:  0.27625
Test Accuracy:  0.55
Test Error:  0.44999999999999996
```

[5.4.6]Confusion Matrix

```
In [110]: from sklearn.metrics import confusion_matrix
conf_matrix = confusion_matrix(Y_Test, optimal_model.predict(tfidf_sent_vector
s_test))
df_conf_matrix = pd.DataFrame(conf_matrix, columns=np.unique(Y_Test), index=np
.unique(Y_Test))
df_conf_matrix.index.name = 'Actual'
df_conf_matrix.columns.name = 'Predicted'
plt.figure(figsize=(10,7))
sns.set(font_scale=1.4)
sns.heatmap(df_conf_matrix, cmap='Blues', annot=True, annot_kws={'size':16}, f
mt='d')
```

Out[110]: <matplotlib.axes._subplots.AxesSubplot at 0x22682768e10>



[5.4.7]Classification Report

```
In [111]: from sklearn.metrics import classification_report
print(classification_report(Y_Test, prediction))
```

	precision	recall	f1-score	support
0	0.53	0.57	0.55	193
1	0.57	0.53	0.55	207
avg / total	0.55	0.55	0.55	400

Pretty Table

```
In [112]: from prettytable import PrettyTable
comparision = PrettyTable()
comparision.field_names = ["Vectorizer", "Algorithm", "Hyperparameter", "AUC"]
comparision.add_row(["BoW", 'brute', optimal_k1, np.round(float(roc_auc1),3)])
comparision.add_row(["TF-IDF", 'brute', optimal_k2, np.round(float(roc_auc3),3)])
comparision.add_row(["Word2Vec", 'brute', optimal_k3, np.round(float(roc_auc5),3)])
comparision.add_row(["TF-IDF Weighted W2V", 'brute', optimal_k4, np.round(float(roc_auc7),3)])
comparision.add_row(["BoW", 'kd-tree', optimal_k1_kdtree, np.round(float(roc_auc9),3)])
comparision.add_row(["TF-IDF", 'kd-tree', optimal_k2_kdtree, np.round(float(roc_auc11),3)])
comparision.add_row(["Word2Vec", 'kd-tree', optimal_k3_kdtree, np.round(float(roc_auc13),3)])
comparision.add_row(["TF-IDF Weighted W2V", 'kd-tree', optimal_k4_kdtree, np.round(float(roc_auc15),3)])
print(comparision)
```

Vectorizer	Algorithm	Hyperparameter	AUC
BoW	brute	19	0.67
TF-IDF	brute	19	0.67
Word2Vec	brute	47	0.5
TF-IDF Weighted W2V	brute	43	0.819
BoW	kd-tree	41	0.739
TF-IDF	kd-tree	41	0.751
Word2Vec	kd-tree	11	0.574
TF-IDF Weighted W2V	kd-tree	5	0.56

Conclusion

Divided the dataset into Train, Test and CV, Used Train dataset to train the each featurization.

Concluded as below from the above pretty table:

Here classifying the query point by considering K as odd number

1. TF-IDF Weighted W2V with brute force model provides best AUC when Hyperparameter is 43.
2. TF-IDF with KD_TREE model provides the best AUC when hypermarameter is 41
3. TF-IDF Weighted W2V AUC score is high as compare to other vectorizer's.