

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 - unique 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 considered as a positive review. A rating of 1 or 2 can be considered as negative one. A review of rating 3 is considered neutral 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 will be set to "positive". Otherwise, it will be set to "negative".

```
In [2]: %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

```

```

C:\Users\user\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWarning: detected Windows; aliasing chunkize to chunkize_serial
  warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")

```

```

In [3]: # 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

```

```

# 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 LIMIT 100000""", 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 -1
    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 (100000, 10)

Out[3]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenomin
2	3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	

```
In [4]: display = pd.read_sql_query("""
SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
FROM Reviews
GROUP BY UserId
HAVING COUNT(*)>1
""", con)
```

```
In [5]: print(display.shape)
display.head()
```

(80668, 7)

Out[5]:

	UserId	ProductId	ProfileName	Time	Score	Text	COUNT(*)
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

	UserId	ProductId	ProfileName	Time	Score	Text	COUNT(*)
4	#oc-R12KPBODL2B5ZD	B007OSBE1U	Christopher P. Presta	1348617600	1	I didnt like this coffee. Instead of telling y...	2

In [6]: `display[display['UserId']=='AZY10LLTJ71NX']`

Out[6]:

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

In [7]: `display['COUNT(*)'].sum()`

Out[7]: 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 [8]: `display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID`

```
""", con)
display.head()
```

Out[8]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenon
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan	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 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 [9]: #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 [10]: #Deduplication of entries
final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time",
"Text"}, keep='first', inplace=False)
final.shape
```

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

```
In [11]: #Checking to see how much % of data still remains
(final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

```
Out[11]: 87.775
```

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 [12]: display= pd.read_sql_query("""
SELECT *
```



```

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

display.head()

```

Out[12]:

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

In [13]: `final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]`

In [14]: *#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()`
(87773, 10)

Out[14]:

```

1    73592
-1   14181
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 observed to be better than Porter Stemming)

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

```
In [15]: # 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)
```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its very hard to find any chicken products made in the USA but they are out there, but this one isnt. Its too bad too bec ause its a good product but I wont take any chances till they know what is going on with the china imports.

=====

The Candy Blocks were a nice visual for the Lego Birthday party but the candy has little taste to it. Very little of the 2 lbs that I bought w ere eaten and I threw the rest away. I would not buy the candy again.

=====

was way to hot for my blood, took a bite and did a jig lol

=====

My dog LOVES these treats. They tend to have a very strong fish oil sme ll. So if you are afraid of the fishy smell, don't get it. But I think my dog likes it because of the smell. These treats are really small in size. They are great for training. You can give your dog several of the se without worrying about him over eating. Amazon's price was much more reasonable than any other retailer. You can buy a 1 pound bag on Amazon for almost the same price as a 6 ounce bag at other retailers. It's def initely worth it to buy a big bag if your dog eats them a lot.

=====

```
In [16]: # 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)
```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its very hard to find any chicken products made in the USA but they are out there, but this one isnt. Its too bad too bec ause its a good product but I wont take any chances till they know what is going on with the china imports.

```
In [17]: # 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)
```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its very hard to find any chicken products made in the USA but they are out there, but this one isnt. Its too bad too because its a good product but I wont take any chances till they know what is going on with the china imports.

=====

The Candy Blocks were a nice visual for the Lego Birthday party but the candy has little taste to it. Very little of the 2 lbs that I bought were eaten and I threw the rest away. I would not buy the candy again.

=====

was way to hot for my blood, took a bite and did a jig lol

=====

My dog LOVES these treats. They tend to have a very strong fish oil smell. So if you are afraid of the fishy smell, don't get it. But I think my dog likes it because of the smell. These treats are really small in size. They are great for training. You can give your dog several of the

se without worrying about him over eating. Amazon's price was much more reasonable than any other retailer. You can buy a 1 pound bag on Amazon for almost the same price as a 6 ounce bag at other retailers. It's definitely worth it to buy a big bag if your dog eats them a lot.

```
In [18]: # 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 [19]: sent_1500 = decontracted(sent_1500)
print(sent_1500)
print("="*50)
```

```
was way to hot for my blood, took a bite and did a jig lol
=====
```

```
In [20]: #remove words with numbers python: https://stackoverflow.com/a/18082370/4084039
sent_0 = re.sub("\S*\d\S*", "", sent_0).strip()
print(sent_0)
```

```
My dogs loves this chicken but its a product from China, so we wont be
buying it anymore. Its very hard to find any chicken products made in
```

the USA but they are out there, but this one isnt. Its too bad too because its a good product but I wont take any chances till they know what is going on with the china imports.

```
In [21]: #remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)
print(sent_1500)
```

was way to hot for my blood took a bite and did a jig lol

```
In [22]: # https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
# <br /><br /> ==> after the above steps, we are getting "br br"
# we are including them into stop words list
# instead of <br /> if we have <br/> these tags would have revmoved in the 1st step

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

```
"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 [23]: # 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%|██████████| 87773/87773 [01:10<00:00, 1238.24it/s]
```

```
In [24]: preprocessed_reviews[1500]
```

```
Out[24]: 'way hot blood took bite jig lol'
```

```
In [25]: final['CleanedText'] = preprocessed_reviews
         final.head(5)
```

Out[25]:

Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
----	-----------	--------	-------------	----------------------	------------------------

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessD
22620	24750	2734888454	A13ISQV0U9GZIC	Sandikaye	1	
22621	24751	2734888454	A1C298ITT645B6	Hugh G. Pritchard	0	
70677	76870	B00002N8SM	A19Q006CSFT011	Arielle	0	
70676	76869	B00002N8SM	A1FYH4S02BW7FN	wonderer	0	
70675	76868	B00002N8SM	AUE8TB5VHS6ZV	eyeofthestorm	0	

[3.2] Preprocessing Review Summary

In [0]: `## Similarly you can do preprocessing for review summary also.`

[4] Featurization

[4.1] BAG OF WORDS

```
In [0]: #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', 'aahhs', 'aback', 'abandon', 'abates', 'abb
ott', 'abby', 'abdominal', 'abiding', 'ability']
=====
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (4986, 12997)
the number of unique words  12997
```

[4.2] Bi-Grams and n-Grams.

```
In [0]: #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 numebtrs min_df=10, max_features=5000, of your ch
oice
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_shape())
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 (4986, 3144)
the number of unique words including both unigrams and bigrams 3144

```

[4.3] TF-IDF

In [0]:

```

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.get_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', 'able find', 'able get', 'absolute', 'absolutely', 'absolutely delicious', 'absolutely love', 'absolutely no', 'according']

```

```

=====

```

```

the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer (4986, 3144)
the number of unique words including both unigrams and bigrams 3144

```

[4.4] Word2Vec

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

```
In [0]: # 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 variable 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 occurred at least 5 times
    w2v_model=Word2Vec(list_of_sentence,min_count=5,size=50, workers=4)
    print(w2v_model.wv.most_similar('great'))
    print('='*50)
    print(w2v_model.wv.most_similar('worst'))

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_train_w2v = True, to train your own w2v ")
```

```
[('snack', 0.9951335191726685), ('calorie', 0.9946465492248535), ('wonderful', 0.9946032166481018), ('excellent', 0.9944332838058472), ('especially', 0.9941144585609436), ('baked', 0.9940600395202637), ('salted', 0.994047224521637), ('alternative', 0.9937226176261902), ('tasty', 0.9936816692352295), ('healthy', 0.9936649799346924)]
=====
[('varieties', 0.9994194507598877), ('become', 0.9992934465408325), ('popcorn', 0.9992750883102417), ('de', 0.9992610216140747), ('miss', 0.9992451071739197), ('melitta', 0.999218761920929), ('choice', 0.9992102384567261), ('american', 0.9991837739944458), ('beef', 0.9991780519485474), ('finish', 0.9991567134857178)]
```

```
In [0]: 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 3817
sample words ['product', 'available', 'course', 'total', 'pretty', 'sticky', 'right', 'nearby', 'used', 'ca', 'not', 'beat', 'great', 'received', 'shipment', 'could', 'hardly', 'wait', 'try', 'love', 'call', 'instead', 'removed', 'easily', 'daughter', 'designed', 'printed', 'use', 'car', 'windows', 'beautifully', 'shop', 'program', 'going', 'lot', 'fun', 'everywhere', 'like', 'tv', 'computer', 'really', 'good', 'idea', 'final', 'outstanding', 'window', 'everybody', 'asks', 'bought', 'made']
```

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

[4.4.1.1] Avg W2v

```
In [0]: # 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_sentence): # 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%|████████████████████████████████████████████████████████████████████████████████| 4986/4986 [00:03<00:00, 1330.47it/s]
```

```
4986
50
```

[4.4.1.2] TFIDF weighted W2v

```
In [0]: # 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_)))
```

[illegible]

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

2. Hyper paramter tuning (find best hyper parameters corresponding the algorithm that you choose)

- 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

3. Pertubation Test

- Get the weights W after fit your model with the data X i.e Train data.
- Add a noise to the X ($X' = X + e$) and get the new data set X' (if X is a sparse matrix, $X.data += e$)
- Fit the model again on data X' and get the weights W'
- Add a small eps value(to eliminate the divisible by zero error) to W and W' i.e $W = W + 10^{-6}$ and $W' = W' + 10^{-6}$
- Now find the % change between W and W' ($| (W - W') / (W) | * 100$)
- Calculate the 0th, 10th, 20th, 30th, ... 100th percentiles, and observe any sudden rise in the values of percentage_change_vector
- Ex: consider your 99th percentile is 1.3 and your 100th percentiles are 34.6, there is sudden rise from 1.3 to 34.6, now calculate the 99.1, 99.2, 99.3,..., 100th percentile values and get the proper value after which there is sudden rise the values, assume it is 2.5
- Print the feature names whose % change is more than a threshold x (in our example it's 2.5)

4. Sparsity

- Calculate sparsity on weight vector obtained after using L1 regularization

NOTE: Do sparsity and multicollinearity for any one of the vectorizers. Bow or tf-idf is

recommended.

5. Feature importance

- Get top 10 important features for both positive and negative classes separately.

6. Feature engineering

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

7. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.



Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.



Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](#).



8. [Conclusion](#)

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



Note: Data Leakage

1. There will be an issue of data-leakage if you vectorize the entire data and then split it into train/cv/test.

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

Applying Logistic Regression

[5.1] Logistic Regression on BOW, SET 1

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

```
In [0]: # Please write all the code with proper documentation
```

```
In [26]: X = final["CleanedText"]
print("shape of X:", X.shape)

shape of X: (87773,)
```

```
In [27]: y = final["Score"]
print("shape of y:", y.shape)

shape of y: (87773,)
```

```
In [28]: 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=False): 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)

print("="*100)

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)

```

```

(39400,) (39400,)
(19407,) (19407,)
(28966,) (28966,)

```

```

=====
=====
After vectorizations
(39400, 37417) (39400,)
(19407, 37417) (19407,)
(28966, 37417) (28966,)
=====
=====

```

```

In [60]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
train_auc = []

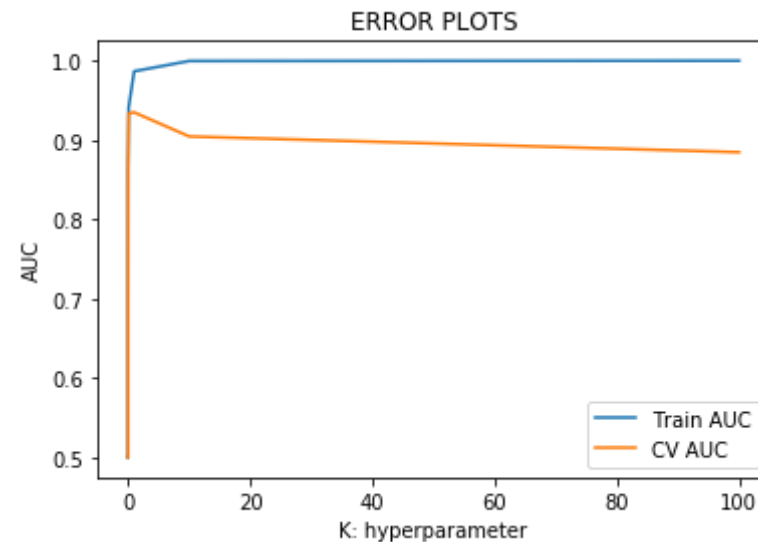
```

```

cv_auc = []
K = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100]
for i in K :
    clf = LogisticRegression(penalty='l1', C=i)
    clf.fit(X_train_bow, y_train)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    y_train_pred = clf.predict_proba(X_train_bow)[:,-1]
    y_cv_pred = clf.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()

```



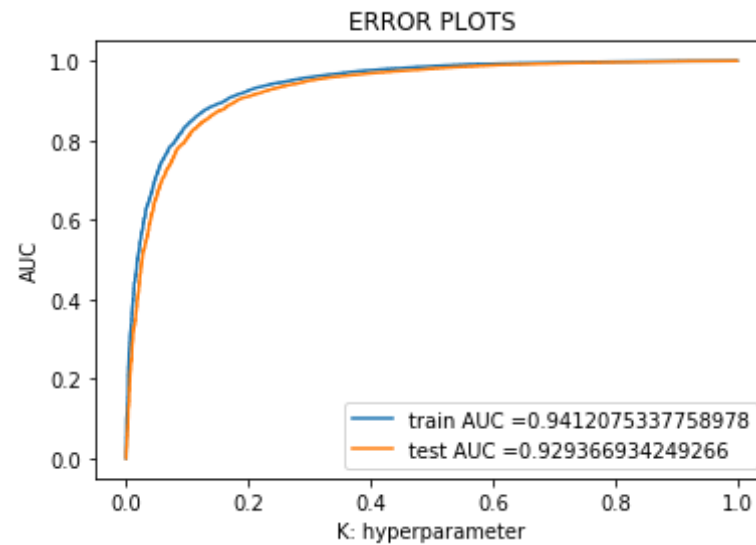
In [61]: best_c=0.1

```
In [62]: from sklearn.metrics import roc_curve, auc

clf = LogisticRegression(penalty='l1',C=best_c)
clf.fit(X_train_bow, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

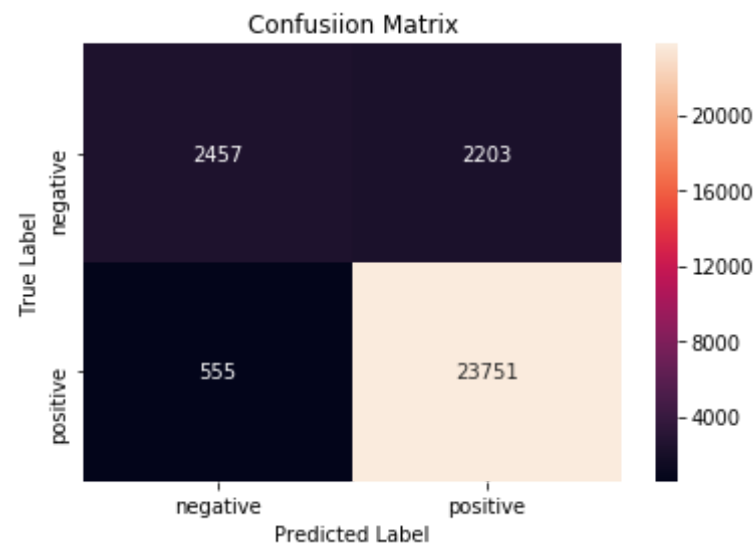
train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba(X_train_bow)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(X_test_bow)[:,1])

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



```
In [63]: print("Test confusion matrix")
cm=confusion_matrix(y_test, clf.predict(X_test_bow))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

Test confusion matrix



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

```
In [0]: # Please write all the code with proper documentation
```

```
In [45]: clf = LogisticRegression(penalty= 'l1',C=100)
         clf.fit(X_train_bow, y_train)
         y_pred = clf.predict_proba(X_test_bow)
         print("Non Zero weights:",np.count_nonzero(clf.coef_))
```

Non Zero weights: 7439

```
In [47]: clf = LogisticRegression(penalty= 'l1',C=10)
         clf.fit(X_train_bow, y_train)
         y_pred = clf.predict_proba(X_test_bow)
         print("Non Zero weights:",np.count_nonzero(clf.coef_))
```

Non Zero weights: 6613

```
In [48]: clf = LogisticRegression(penalty= 'l1',C=1)
```

```
clf.fit(X_train_bow, y_train)
y_pred = clf.predict_proba(X_test_bow)
print("Non Zero weights:", np.count_nonzero(clf.coef_))
```

Non Zero weights: 3486

```
In [50]: clf = LogisticRegression(penalty= 'l1', C=.1)
         clf.fit(X_train_bow, y_train)
         y_pred = clf.predict_proba(X_test_bow)
         print("Non Zero weights:", np.count_nonzero(clf.coef_))
```

Non Zero weights: 657

```
In [51]: clf = LogisticRegression(penalty= 'l1', C=.01)
         clf.fit(X_train_bow, y_train)
         y_pred = clf.predict_proba(X_test_bow)
         print("Non Zero weights:", np.count_nonzero(clf.coef_))
```

Non Zero weights: 80

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

```
In [0]: # Please write all the code with proper documentation
```

```
In [88]: 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=False): 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)

print("="*100)

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)

```

```

(39400,) (39400,)
(19407,) (19407,)
(28966,) (28966,)

```

```

=====
=====
After vectorizations
(39400, 37632) (39400,)
(19407, 37632) (19407,)
(28966, 37632) (28966,)
=====
=====

```

```

In [53]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
train_auc = []
cv_auc = []
K = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10,100]
for i in K :
    clf = LogisticRegression(penalty='l2',C=i)

```

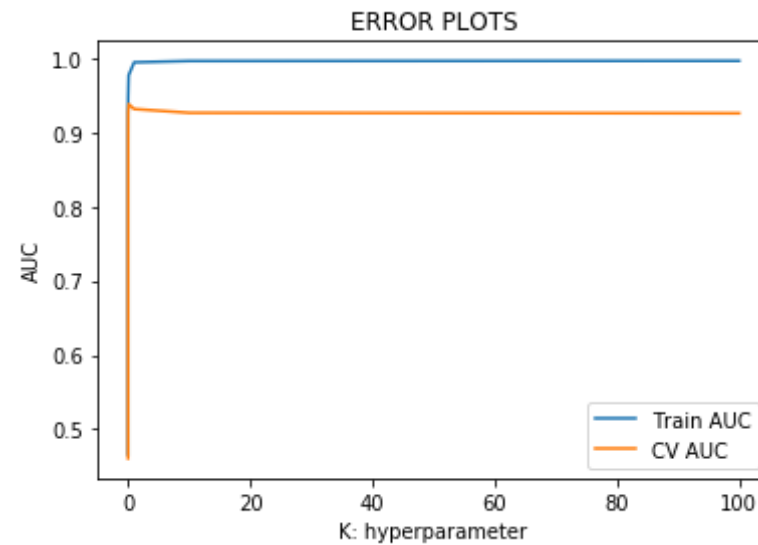


```

clf.fit(X_train_bow, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs
y_train_pred = clf.predict_proba(X_train_bow)[:,-1]
y_cv_pred = clf.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 [54]: `best_c=0.1`

In [55]: `from sklearn.metrics import roc_curve, auc`

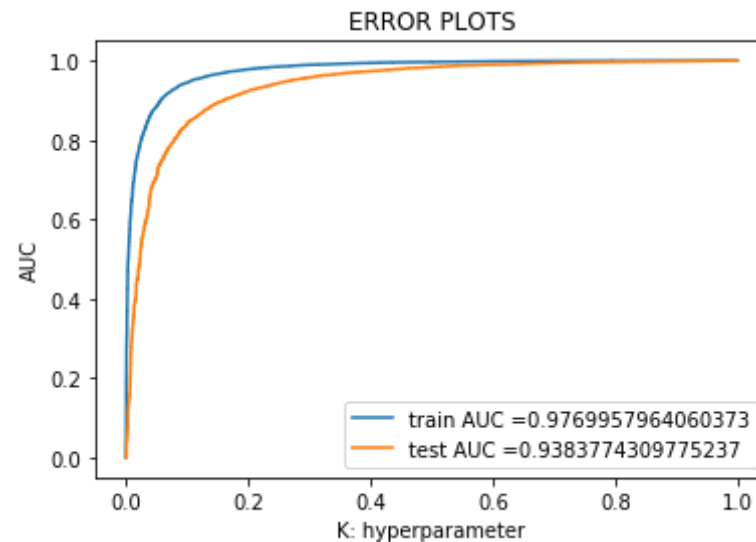
```

clf = LogisticRegression(penalty='l2',C=best_c)
clf.fit(X_train_bow, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

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

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

```



[5.1.2.1] Performing perturbation test (multicollinearity check) on BOW, SET 1

```
In [0]: # Please write all the code with proper documentation
```

```
In [29]: from sklearn.linear_model import LogisticRegression
clf = LogisticRegression(penalty='l2',C=10)
clf.fit(X_train_bow, y_train)
y_pred = clf.predict_proba(X_test_bow)
```

```
In [30]: from scipy.sparse import find
#Weights before adding random noise
weights1 = find(clf.coef_[0])[2]
print(weights1.shape)
print(weights1[:50])
```

```
(37417,)
[ 8.74926546e-01  8.25829399e-02  1.32284493e-01  1.50364184e-03
 1.00069805e-05  7.70190630e-05  7.70190630e-05  1.57473484e-05
 4.58572336e-06  5.74335573e-01  7.87528095e-05  7.34965241e-06
 1.79862745e-04 -4.91623935e-05  4.03433258e-03 -5.38618371e-01
 1.91250457e-02  8.49451789e-01  3.41170719e-02 -2.27616651e+00
 2.38509353e-02  2.56791861e-03  1.62627434e-05  2.41194168e-03
-1.45539670e-08  7.80746059e-02  4.44583762e-01 -9.34171177e-01
 2.86147491e-02 -7.35281888e-01 -1.00209271e+00 -8.22679476e-02
 3.66120289e-03  7.03937100e-05  1.39152806e-01 -3.31495537e-02
 3.31781731e-01  4.04114893e-01 -7.98643245e-01  5.13865172e-04
 9.96359915e-01 -7.01509369e-01 -2.30254724e-02 -4.30519637e-01
 1.10882523e-03  8.24539505e-02 -8.22679476e-02  6.50835103e-05
-2.09458493e-01 -9.46926439e-04]
```

```
In [31]: X_train_t = X_train_bow
#Random noise
epsilon = 0.0001
#Getting the postions(row and column) and value of non-zero datapoints

#Introducing random noise to non-zero datapoints
X_train_t.data = epsilon + X_train_t.data
```

```
In [32]: from sklearn.linear_model import LogisticRegression
```

```
cl = LogisticRegression(penalty= 'l2',C=10)
cl.fit(X_train_t,y_train)
y_pred = cl.predict_proba(X_test_bow)
```

```
In [33]: weights2 = find(cl.coef_[0])[2]
```

```
print(weights2.shape)
print(weights2[:50])
```

```
(37417,)
[ 8.75081728e-01  8.15160647e-02  1.32145382e-01  1.50243024e-03
 1.58724684e-05  1.22734847e-04  1.22734847e-04  2.29009333e-05
 7.58135662e-06  5.75090044e-01  8.11589850e-05  1.22025408e-05
 2.72099371e-04 -6.39952625e-05  4.12651003e-03 -5.40088472e-01
 1.91210385e-02  8.50513232e-01  3.43170881e-02 -2.27328376e+00
 2.38168515e-02  2.67031139e-03  2.69995580e-05  2.09032598e-03
-2.43062145e-08  7.82135021e-02  4.42501284e-01 -9.34164028e-01
 2.90389778e-02 -7.33950140e-01 -1.00165734e+00 -8.24433245e-02
 3.67900877e-03  9.78273875e-05  1.39499160e-01 -3.28288333e-02
 3.32552503e-01  4.02344858e-01 -7.99333577e-01  8.53383807e-04
 9.96466557e-01 -7.01561556e-01 -2.31734659e-02 -4.30581374e-01
 1.19978083e-03  8.21192447e-02 -8.24433245e-02  8.35876072e-05
-2.08833527e-01 -9.62887362e-04]
```

```
In [34]: weights1=weights1+.000001
weights2=weights2+.000001
```

```
In [35]: weights_diff =abs((weights1 - weights2)/weights1) * 100
```

```
In [36]: for i in range(10, 101, 10):
          print("{}th Percentile value : {:.5f}".format(i, np.percentile(weights_diff, i)))
```

```
10th Percentile value : 0.03293
20th Percentile value : 0.07151
30th Percentile value : 0.12547
40th Percentile value : 0.22374
```

```
50th Percentile value : 0.42073
60th Percentile value : 0.99445
70th Percentile value : 3.60170
80th Percentile value : 21.12104
90th Percentile value : 54.37340
100th Percentile value : 5621.91804
```

There is a sudden raise from 90th percentile to 100th percentile

```
In [37]: for i in range(90, 101, 1):
          print("{}th Percentile value : {:.5f}".format(i, np.percentile(weights_diff, i)))
```

```
90th Percentile value : 54.37340
91th Percentile value : 56.57563
92th Percentile value : 58.58371
93th Percentile value : 60.26458
94th Percentile value : 62.06932
95th Percentile value : 63.44236
96th Percentile value : 64.93060
97th Percentile value : 66.37869
98th Percentile value : 74.66525
99th Percentile value : 122.53356
100th Percentile value : 5621.91804
```

There is a sudden raise from 99th percentile to 100th percentile

```
In [53]: for i in np.linspace(99,100,11).tolist():
          print("{}th Percentile value : {:.5f}".format(i, np.percentile(weights_diff, i)))
```

```
99.0th Percentile value : 122.53356
99.1th Percentile value : 134.55448
99.2th Percentile value : 145.05258
99.3th Percentile value : 154.07688
99.4th Percentile value : 173.56011
99.5th Percentile value : 188.97404
99.6th Percentile value : 223.64101
```

```
99.7th Percentile value : 294.32792
99.8th Percentile value : 379.57960
99.9th Percentile value : 739.95890
100.0th Percentile value : 5621.91804
```

There is a sudden raise from 99.9th percentile to 100th percentile

```
In [56]: percentchange = pd.DataFrame(weights_diff, index = vectorizer.get_feature_names(), columns=['Change'])
percentchange.head(5)
```

Out[56]:

	Change
aa	0.017737
aaa	1.291868
aaaa	0.105160
aaaaa	0.080524
aaaaaaaaaaa	53.288801

```
In [58]: percentchange = percentchange[percentchange['Change'] > 739]
percentchange.shape
```

Out[58]: (39, 1)

```
In [67]: percentchange.sort_values(by='Change', ascending=False)
percentchange.head()
```

Out[67]:

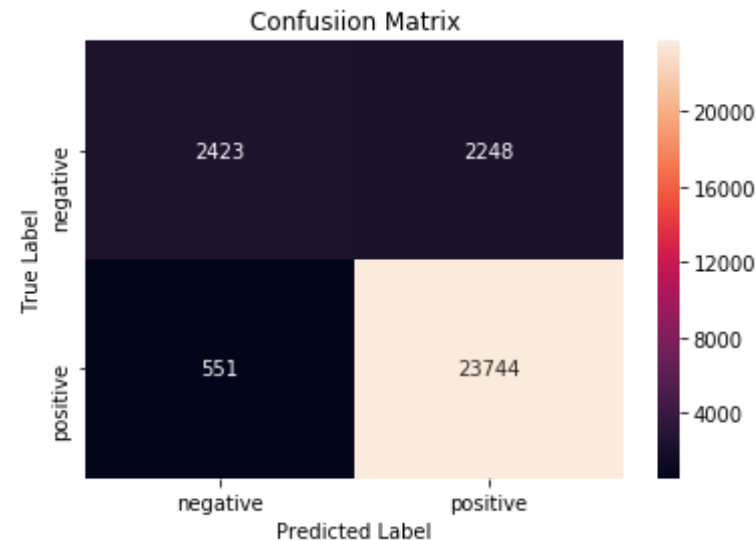
	Change
glase	5621.918037
stratch	5621.918037
dolce	4706.814569
punctuated	4022.067357

Change

absolutly 3590.639699

```
In [127]: print("Test confusion matrix")
cm=confusion_matrix(y_test, clf.predict(X_test_bow))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

Test confusion matrix



[5.1.3] Feature Importance on BOW, SET 1

[5.1.3.1] Top 10 important features of positive class from SET 1

```
In [0]: # Please write all the code with proper documentation
```

```
In [126]: feature_names = vectorizer.get_feature_names()
coefs_with_fns = sorted(zip(clf.coef_[0], feature_names))
top = (coefs_with_fns[:-(10 + 1):-1])
print("\tPositive")
for (coef_2, fn_2) in top:
    print("\t%.4f\t%-15s\t\t\t\t" % (coef_2, fn_2))
```

```
Positive
1.4194 delicious
1.3813 amazing
1.3709 perfect
1.2808 excellent
1.2205 yummy
1.1592 wonderful
1.1399 loves
1.1153 highly
1.0842 pleased
1.0723 best
```

[5.1.3.2] Top 10 important features of negative class from SET 1

```
In [0]: # Please write all the code with proper documentation
```

```
In [125]: feature_names = vectorizer.get_feature_names()
coefs_with_fns = sorted(zip(clf.coef_[0], feature_names))
top = (coefs_with_fns[:10])
print("\tNegative")

for (coef_1, fn_1) in top:
    print("\t%.4f\t%-15s\t\t\t\t" % (coef_1, fn_1))
```

```
Negative
-2.3314 worst
```



```
-1.8924 terrible
-1.8326 disappointing
-1.6010 awful
-1.3771 threw
-1.3659 yuck
-1.3407 horrible
-1.3027 disgusting
-1.2750 disappointment
-1.2620 disappointed
```

[5.2] Logistic Regression on TFIDF, SET 2

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

```
In [114]: # Please write all the code with proper documentation
```

```
In [64]: X = final["CleanedText"]
print("shape of X:", X.shape)

shape of X: (87773,)
```

```
In [65]: y = final["Score"]
print("shape of y:", y.shape)

shape of y: (87773,)
```

```
In [66]: 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=False): 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)

print("="*100)

from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(ngram_range=(1,2))
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_tfidf = vectorizer.transform(X_train)
X_cv_tfidf = vectorizer.transform(X_cv)
X_test_tfidf = vectorizer.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)

```

```

(39400,) (39400,)
(19407,) (19407,)
(28966,) (28966,)

```

```

=====
=====
After vectorizations
(39400, 773213) (39400,)
(19407, 773213) (19407,)
(28966, 773213) (28966,)
=====
=====

```

```

In [67]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
train_auc = []

```

```

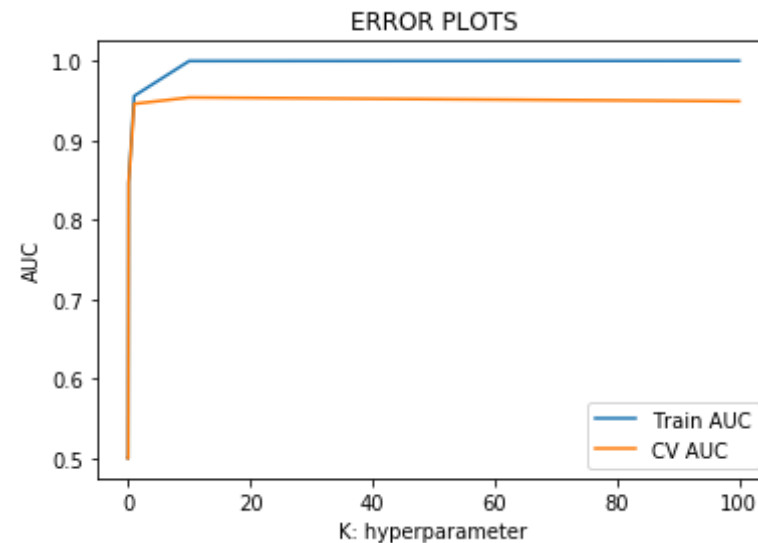
cv_auc = []

K = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100]
for i in K :
    clf=LogisticRegression(penalty='l1',C=i)
    clf.fit(X_train_tfidf, y_train)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    y_train_pred = clf.predict_proba(X_train_tfidf)[:,-1]
    y_cv_pred = clf.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()

```



```
In [68]: best_c=1
```

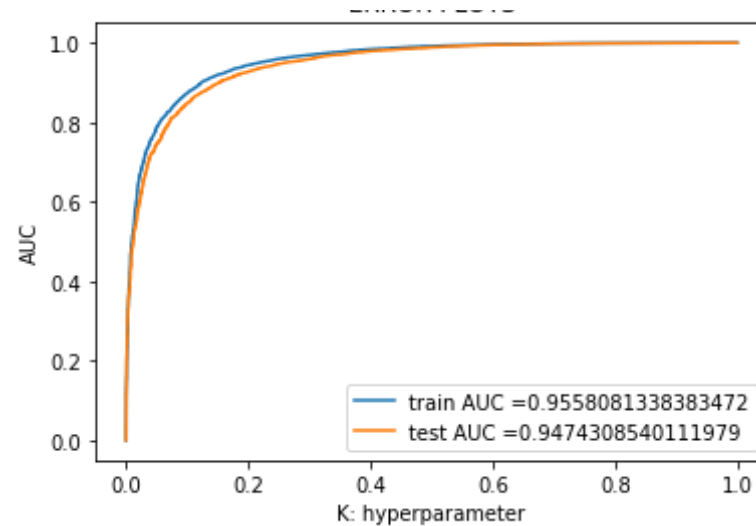
```
In [69]: from sklearn.metrics import roc_curve, auc

clf=LogisticRegression(penalty='l1',C=best_c)
clf.fit(X_train_tfidf, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba(X_train_tfidf)[:,-1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(X_test_tfidf)[:,-1])

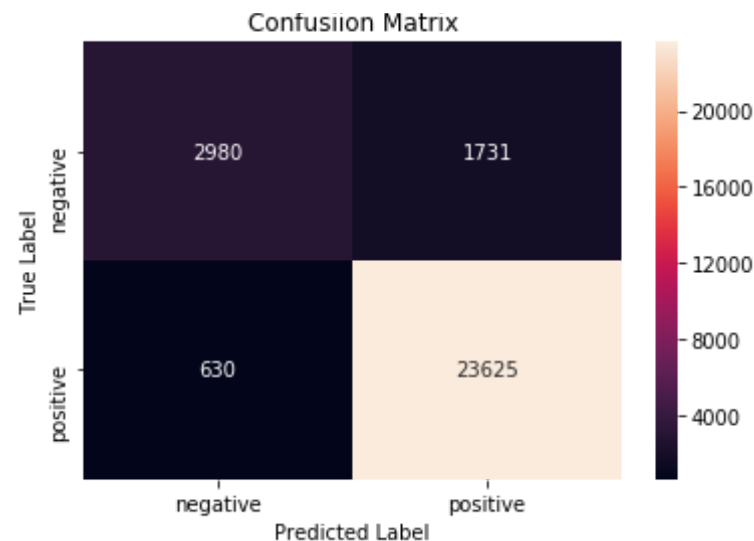
plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```

ERROR PLOTS



```
In [70]: from sklearn.metrics import confusion_matrix
import seaborn as sns
print("Test confusion matrix")
cm=confusion_matrix(y_test, clf.predict(X_test_tfidf))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

Test confusion matrix



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

In [115]: *# Please write all the code with proper documentation*

```
In [71]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
train_auc = []
cv_auc = []

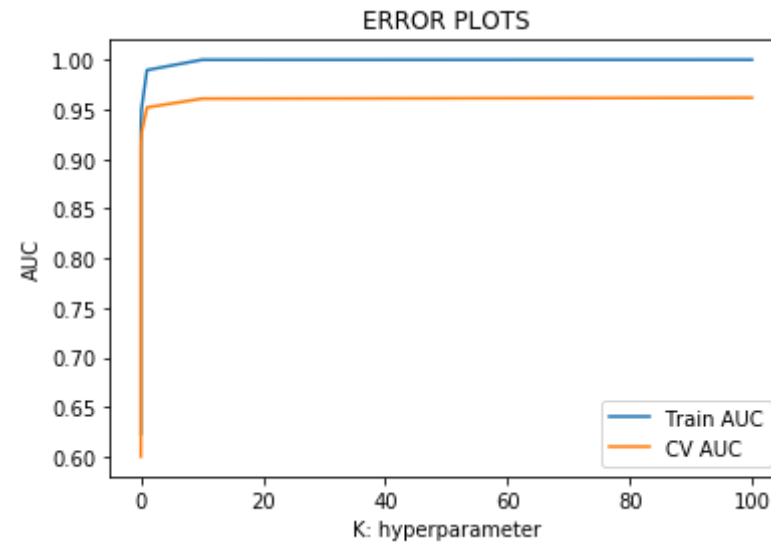
K = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100]
for i in K :
    clf=LogisticRegression(penalty='l2',C=i)
    clf.fit(X_train_tfidf, y_train)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    y_train_pred = clf.predict_proba(X_train_tfidf)[:,-1]
    y_cv_pred = clf.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()

```



In [72]: `best_c=1`

```

In [73]: from sklearn.metrics import roc_curve, auc

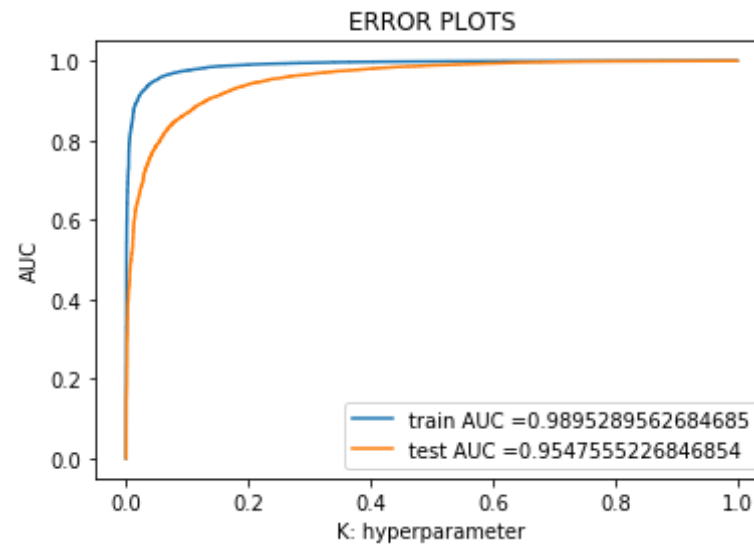
clf=LogisticRegression(penalty='l2',C=best_c)
clf.fit(X_train_tfidf, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability
  estimates of the positive class

```

```
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba(
X_train_tfidf)[: ,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(X_
test_tfidf)[: ,1])

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

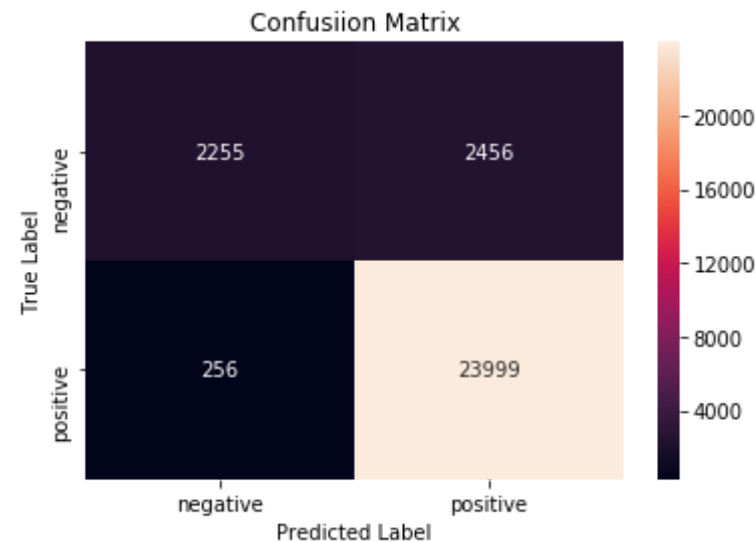


```
In [74]: from sklearn.metrics import confusion_matrix
import seaborn as sns
print("Test confusion matrix")
cm=confusion_matrix(y_test, clf.predict(X_test_tfidf))
class_label = ["negative", "positive"]
```



```
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

Test confusion matrix



[5.2.3] Feature Importance on TFIDF, SET 2

[5.2.3.1] Top 10 important features of positive class from SET 2

```
In [116]: # Please write all the code with proper documentation
```

```
In [143]: feature_names = vectorizer.get_feature_names()
coefs_with_fns = sorted(zip(clf.coef_[0], feature_names))
top = (coefs_with_fns[:-(10 + 1):-1])
print("\tPositive")
```

```
for (coef_2, fn_2) in top:
    print("\t%.4f\t%-15s\t\t\t\t" % (coef_2, fn_2))
```

```
Positive
11.6879 great
7.3871 good
7.3512 best
7.2342 delicious
6.7641 love
5.9315 perfect
5.8559 loves
5.0774 nice
4.8873 favorite
4.5853 wonderful
```

[5.2.3.2] Top 10 important features of negative class from SET 2

In [117]: *# Please write all the code with proper documentation*

```
In [144]: feature_names = vectorizer.get_feature_names()
coefs_with_fns = sorted(zip(clf.coef_[0], feature_names))
top = (coefs_with_fns[:10])
print("\tNegative")

for (coef_1, fn_1) in top:
    print("\t%.4f\t%-15s\t\t\t\t" % (coef_1, fn_1))
```

```
Negative
-10.2222 not
-7.7723 disappointed
-5.7578 worst
-5.3640 terrible
-5.3554 bad
-5.0556 awful
-5.0188 not good
-4.9472 money
-4.9395 not buy
-4.7779 stale
```

[5.3] Logistic Regression on AVG W2V, SET 3

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

```
In [0]: # Please write all the code with proper documentation
```

```
In [25]: i=0
list_of_sentence=[]
for sentence in final['CleanedText']:
    list_of_sentence.append(sentence.split())
```

```
In [26]: 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 occurred at least 5 times
    w2v_model=Word2Vec(list_of_sentence,min_count=5,size=50, workers=4)
    print(w2v_model.wv.most_similar('great'))
    print('='*50)
    print(w2v_model.wv.most_similar('worst'))

[('awesome', 0.8405764698982239), ('good', 0.8389686346054077), ('fantastic', 0.8386562466621399), ('excellent', 0.818621039390564), ('terrific', 0.8073654770851135), ('wonderful', 0.7965986132621765), ('perfect', 0.7332897782325745), ('amazing', 0.7298626899719238), ('nice', 0.7179173231124878), ('decent', 0.6870323419570923)]
=====
[('greatest', 0.8212548494338989), ('tastiest', 0.7431558966636658), ('best', 0.7354905605316162), ('nastiest', 0.7142645120620728), ('disgusting', 0.6549212336540222), ('horrible', 0.6387206315994263), ('terrible', 0.6338681578636169), ('closest', 0.6336538791656494), ('awful', 0.6313574910163879), ('vile', 0.6184582114219666)]
```



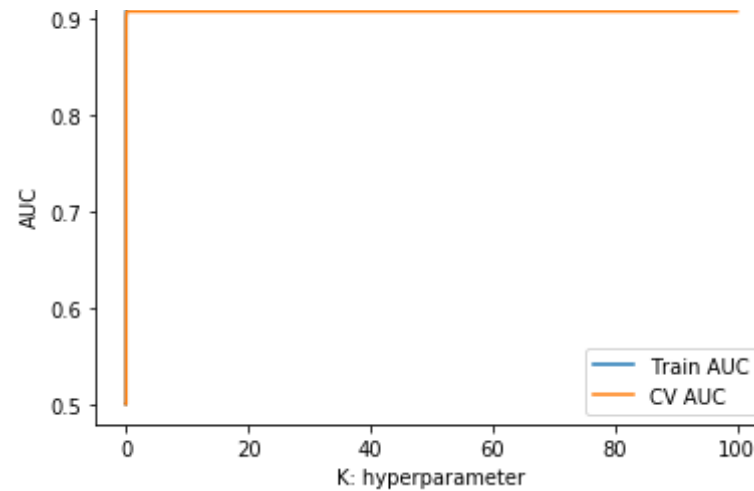
```
# X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=
0.33, shuffle=False): this is for time series split
X_train, X_test, y_train, y_test = train_test_split(sent_vectors, final
['Score'], test_size=0.33) # this is random splitting
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_
size=0.33) # this is random splitting
```

```
In [30]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
train_auc = []
cv_auc = []
K = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100]
for i in K:
    clf=LogisticRegression(penalty='l1',C=i)
    clf.fit(X_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 = clf.predict_proba(X_train)[:,-1]
    y_cv_pred = clf.predict_proba(X_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.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 [31]: `best_c=0.1`

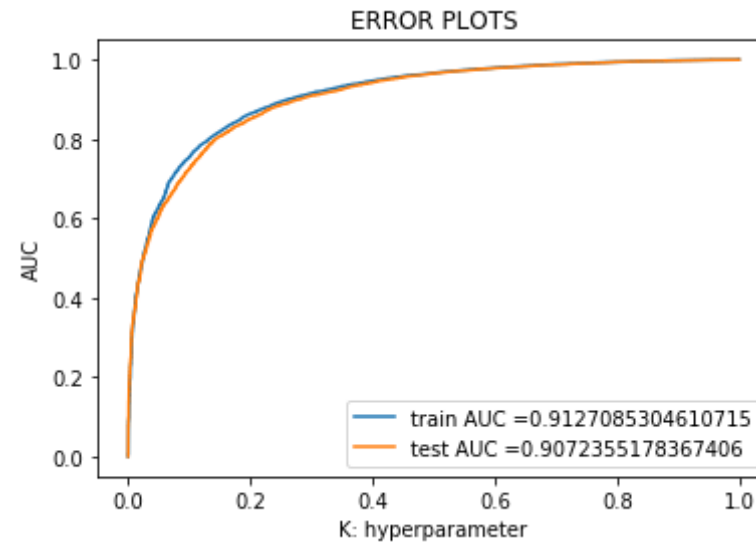
```
In [32]: from sklearn.metrics import roc_curve, auc

clf=LogisticRegression(penalty='l1',C=best_c)
clf.fit(X_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability
# estimates of the positive class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba(
X_train)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(X_
test)[:,1])

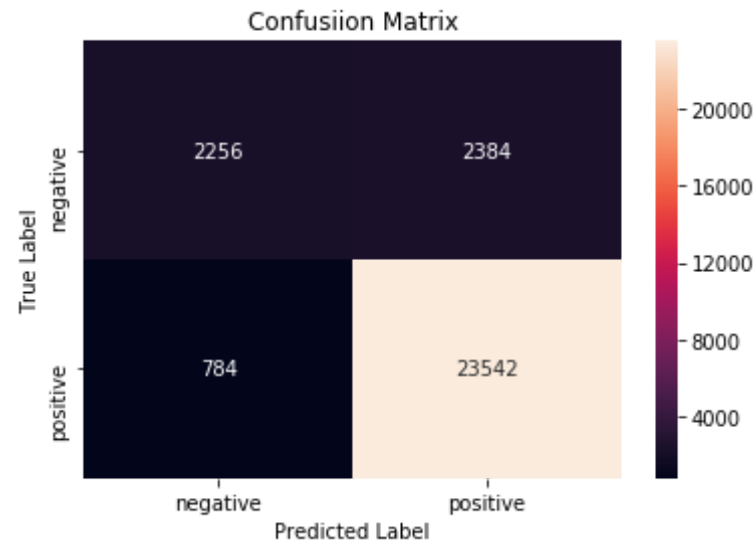
plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, t
rain_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_
tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
```

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



```
In [34]: from sklearn.metrics import confusion_matrix
import seaborn as sns
print("Test confusion matrix")
cm=confusion_matrix(y_test, clf.predict(X_test))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

Test confusion matrix



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

In [0]: *# Please write all the code with proper documentation*

```
In [35]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
train_auc = []
cv_auc = []
K = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100]
for i in K:
    clf=LogisticRegression(penalty='l2',C=i)
    clf.fit(X_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 = clf.predict_proba(X_train)[:,:1]
    y_cv_pred = clf.predict_proba(X_cv)[:,:1]
```

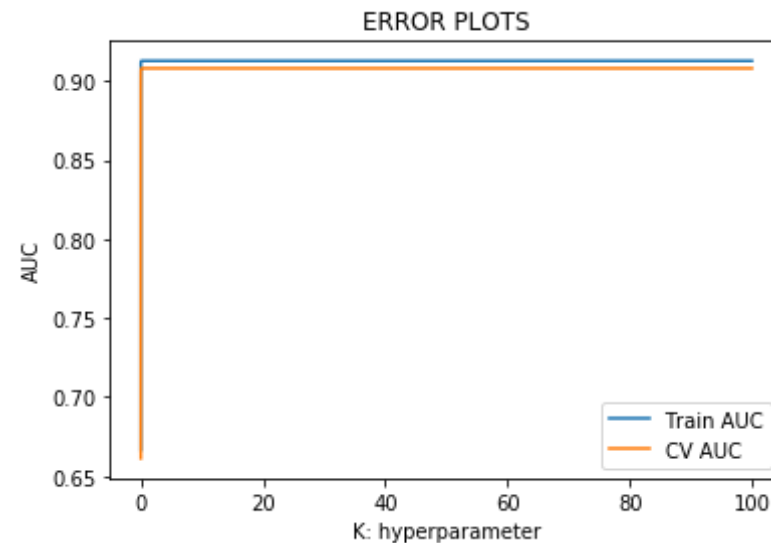


```

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

```



In [36]: `best_c=.1`

In [38]: `from sklearn.metrics import roc_curve, auc`

```

clf=LogisticRegression(penalty='l2',C=best_c)
clf.fit(X_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

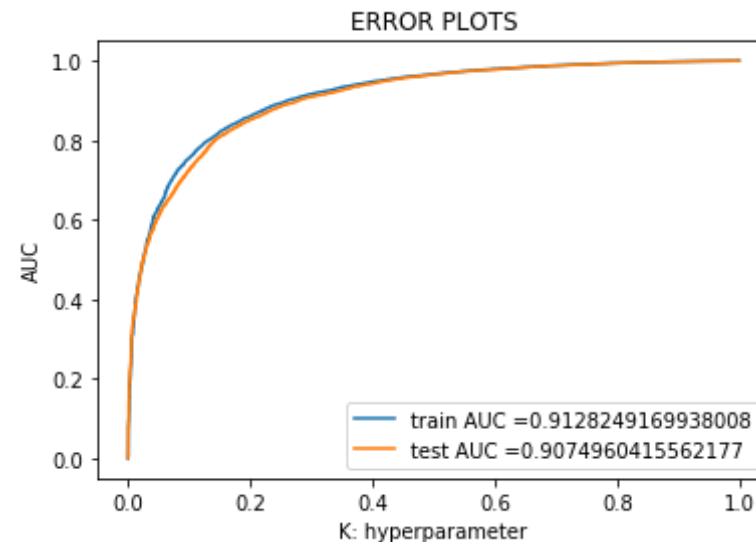
```

```

train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba(X_train)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(X_test)[:,1])

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

```



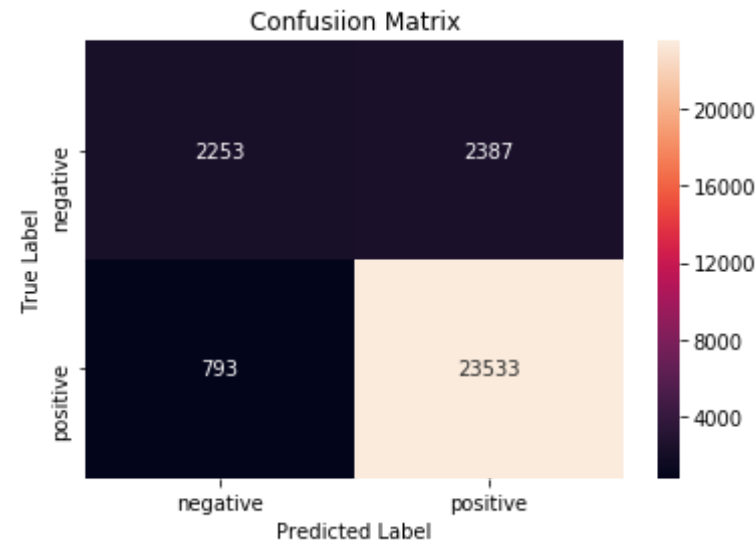
```

In [39]: from sklearn.metrics import confusion_matrix
import seaborn as sns
print("Test confusion matrix")
cm=confusion_matrix(y_test, clf.predict(X_test))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)

```

```
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

Test confusion matrix



[5.4] Logistic Regression on TFIDF W2V, SET 4

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

In [0]: *# Please write all the code with proper documentation*

In [40]: *# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
 model = TfidfVectorizer()
 tf_idf_matrix = model.fit_transform(final['CleanedText'])
 # we are converting a dictionary with word as a key, and the idf as a v*

```
dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
```

```
In [41]: 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_sentence): # 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
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.append(sent_vec)
    row += 1
```

```
100%|███████████|  
87773/87773 [1:35:30<00:00, 15.32it/s]
```

```
In [42]: from sklearn.model_selection import train_test_split

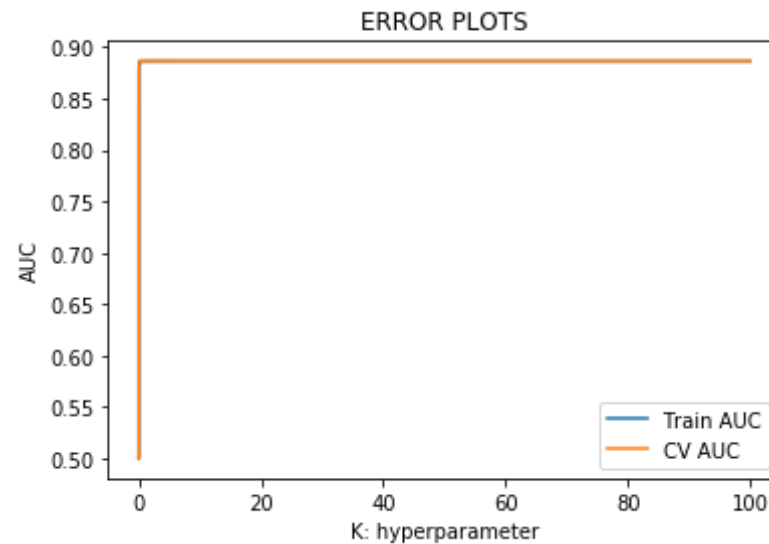
# X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=
0.33, shuffle=False): this is for time series split
X_train, X_test, y_train, y_test = train_test_split(tfidf_sent_vectors,
final['Score'], test_size=0.33) # this is random splitting
```

```
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33)
```

```
In [43]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
train_auc = []
cv_auc = []
K = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100]
for i in K:
    clf=LogisticRegression(penalty='l1',C=i)
    clf.fit(X_train, y_train)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    y_train_pred = clf.predict_proba(X_train)[:,-1]
    y_cv_pred = clf.predict_proba(X_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.plot(K, cv_auc, label='CV AUC')
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```



In [44]: `best_c=0.1`

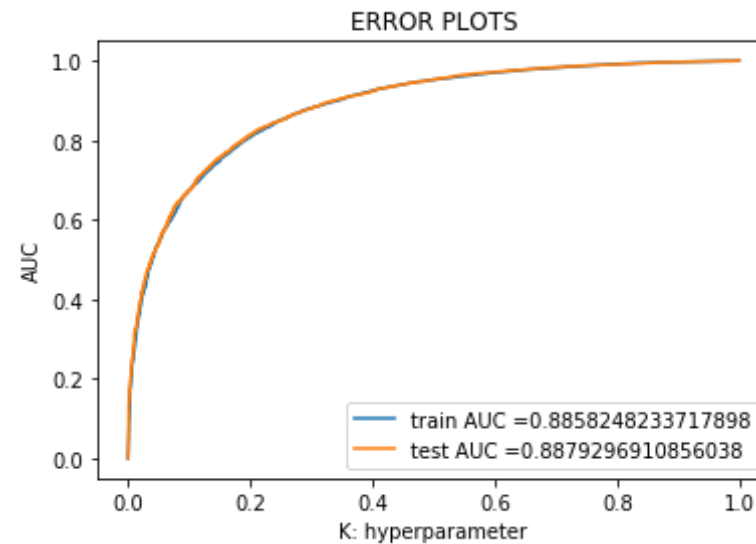
```
In [45]: from sklearn.metrics import roc_curve, auc

clf=LogisticRegression(penalty='l1',C=best_c)
clf.fit(X_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability
# estimates of the positive class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba(
X_train)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(X_
test)[:,1])

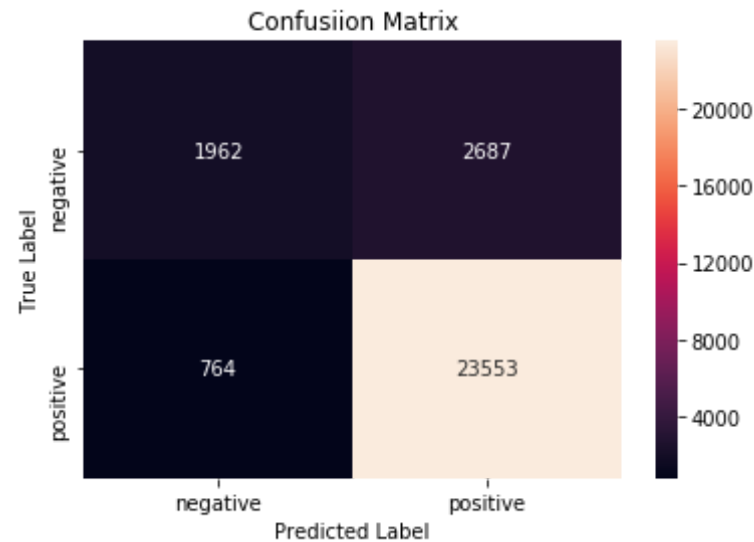
plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, t
rain_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_
tpr)))
plt.legend()
```

```
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```



```
In [46]: from sklearn.metrics import confusion_matrix
import seaborn as sns
print("Test confusion matrix")
cm=confusion_matrix(y_test, clf.predict(X_test))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

Test confusion matrix



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

```
In [0]: # Please write all the code with proper documentation
```

```
In [47]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
train_auc = []
cv_auc = []
K = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100]
for i in K:
    clf=LogisticRegression(penalty='l2',C=i)
    clf.fit(X_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 = clf.predict_proba(X_train)[:,:1]
    y_cv_pred = clf.predict_proba(X_cv)[:,:1]
```

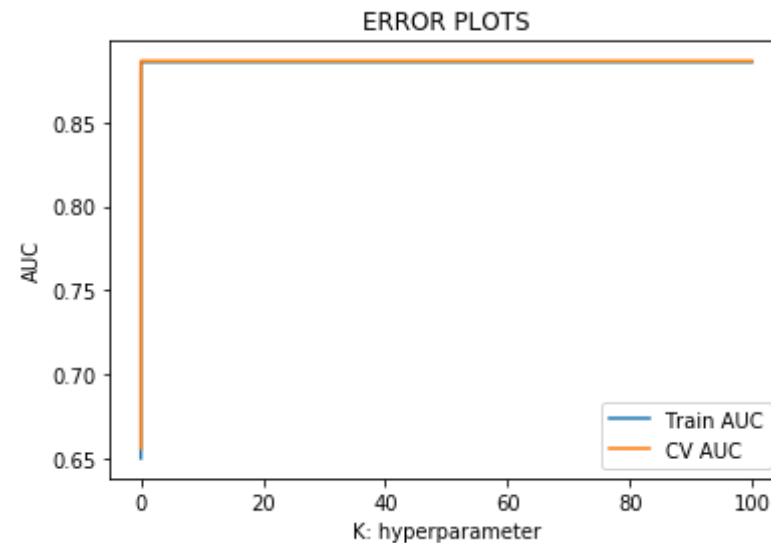


```

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

```



In [48]: `best_c=0.1`

In [49]: `from sklearn.metrics import roc_curve, auc`

```

clf=LogisticRegression(penalty='l2',C=best_c)
clf.fit(X_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

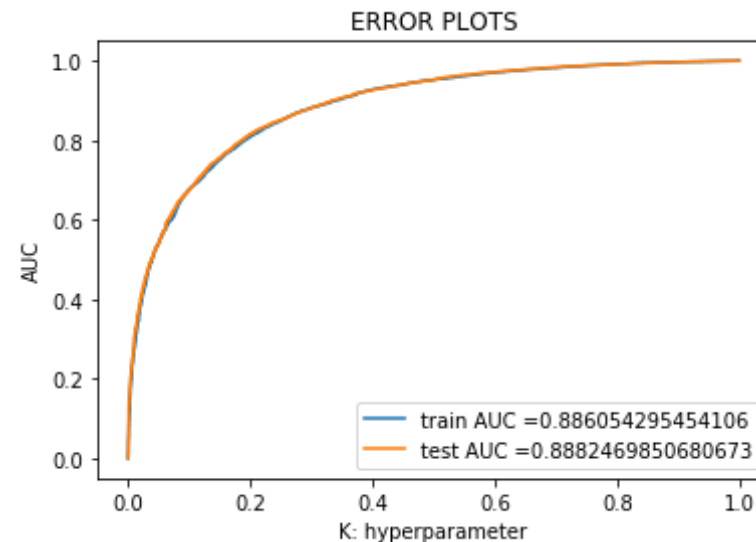
```

```

train_fpr, train_tpr, thresholds = roc_curve(y_train, clf.predict_proba(X_train)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, clf.predict_proba(X_test)[:,1])

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

```



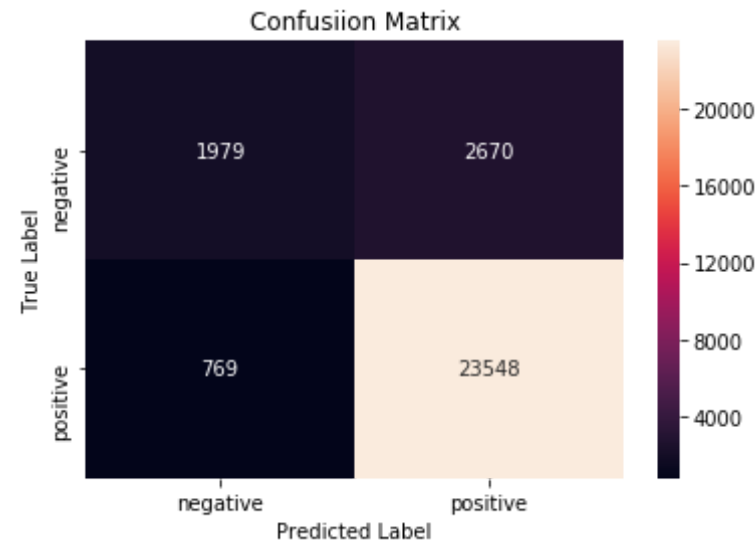
```

In [50]: from sklearn.metrics import confusion_matrix
import seaborn as sns
print("Test confusion matrix")
cm=confusion_matrix(y_test, clf.predict(X_test))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)

```

```
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

Test confusion matrix



[6] Conclusions

In [0]: *# Please compare all your models using Prettytable library*

```
In [54]: models = pd.DataFrame({'vectorizer': ['logisticregression with Bow', "l  
ogisticregression with TFIDF", "logisticregression with Avg_w2v", "logi  
sticregression with tfidf_w2v"], 'Regulatization' : ["l1","l1","l1","l  
1"], 'Hyper Parameter(lamda)': [0.1,1,0.1,0.1], 'AUC': [.93,.94,.90,.88  
]}, columns = ["vectorizer", "Regulatization", "Hyper Parameter(lamda)",  
"AUC"])
models
```

Out[54]:

	vectorizer	Regulatization	Hyper Parameter(lamda)	AUC
0	logisticregression with Bow	l1	0.1	0.93
1	logisticregression with TFIDF	l1	1.0	0.94
2	logisticregression with Avg_w2v	l1	0.1	0.90
3	logisticregression with tfidf_w2v	l1	0.1	0.88

```
In [56]: models = pd.DataFrame({'vectorizer': ['logisticregression with Bow', "logisticregression with TFIDF", "logisticregression with Avg_w2v", "logisticregression with tfidf_w2v"], 'Regulatization': ["l2", "l2", "l2", "l2"], 'Hyper Parameter(lamda)': [0.1, 1, 0.1, 0.1], 'AUC': [.93, .95, .90, .88]}, columns = ["vectorizer", "Regulatization", "Hyper Parameter(lamda)", "AUC"])
models
```

Out[56]:

	vectorizer	Regulatization	Hyper Parameter(lamda)	AUC
0	logisticregression with Bow	l2	0.1	0.93
1	logisticregression with TFIDF	l2	1.0	0.95
2	logisticregression with Avg_w2v	l2	0.1	0.90
3	logisticregression with tfidf_w2v	l2	0.1	0.88

In []: