

# 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:

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 [370]: %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
```

```
In [371]: # 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 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 LIMIT 500000 """, 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 (100000, 10)

Out[371]:

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

```
In [373]: print(display.shape)
display.head()
```

```
(80668, 7)
```

```
Out[373]:
```

	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-R11D9D7SHXIJ9	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 [374]: display[display['UserId']=='AZY10LLTJ71NX']
```

```
Out[374]:
```

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

```
Out[375]: 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 [376]: display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID
""", con)
display.head()
```

Out[376]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenomir
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 [377]: #Sorting data according to ProductId in ascending order
sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True, inplace=False)
```

```
In [378]: #Deduplication of entries
final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"}, keep='first', inplace=True)
final.shape
```

```
Out[378]: (87775, 10)
```

```
In [379]: #Checking to see how much % of data still remains
(final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

```
Out[379]: 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 [380]: display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND Id=44737 OR Id=64422
ORDER BY ProductID
""", con)

display.head()
```

```
Out[380]:
```

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

0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	
---	-------	------------	----------------	-------------------------------	---	--

1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	
---	-------	------------	----------------	-----	---	--



```
In [381]: final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [382]: #Before starting the next phase of preprocessing Lets see the number of entries Le
print(final.shape)

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

(87773, 10)
```

```
Out[382]: 1    73592
          0    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 [383]: # 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 because its a good product but I wont take any chances till they know what is going on with the china imports.

=====

```
In [384]: # 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)
```

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 [385]: # https://stackoverflow.com/questions/16206380/python-beautifulsoup-how-to-remove-
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.

=====

```
In [386]: # 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 [387]: # sent_1500 = decontracted(sent_1500)
# print(sent_1500)
# print("="*50)
```

```
In [388]: #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 [389]: # #remove spacial character: https://stackoverflow.com/a/5843547/4084039
# sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)
# print(sent_1500)
```

```
In [390]: # 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', 'ours',
                "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'I',
                'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itsel',
                'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that',
                'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has',
                'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because',
                'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'th',
                'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off',
                'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all',
                'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than',
                's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've",
                've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "dic",
                "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma',
                "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't",
                'won', "won't", 'wouldn', "wouldn't"])
```

```
In [391]: # Combining all the above stundents
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 [00:43<00:00, 2008.02it/s]

```
In [392]: # preprocessed_reviews[1500]
```

## [3.2] Preprocessing Review Summary

In [393]: *## Similarly you can do preprocessing for review summary also.*

```
In [394]: from sklearn.model_selection import train_test_split

print(type(final['Text']))
final['Text'] = preprocessed_reviews

# Added new feature TextLength in our preprocessed data
final['TextLength'] = final['Text'].str.len()

X = final['Text'].values
Y = final['Score'].values

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.3)

print(X_train.shape, Y_train.shape)
print(X_test.shape, Y_test.shape)

<class 'pandas.core.series.Series'>
(61441,) (61441,)
(26332,) (26332,)
```

## [4] Featurization

### [4.1] BAG OF WORDS

```
In [395]: #Bow
count_vect = CountVectorizer() #in scikit-learn
count_vect.fit(X_train)
print("some feature names ", count_vect.get_feature_names()[:10])
print('='*50)

X_train_bow = count_vect.transform(X_train)
X_test_bow = count_vect.transform(X_test)

print("the type of count vectorizer ", type(X_train_bow))
print("the shape of out text BOW vectorizer ", X_train_bow.get_shape())
print("the number of unique words ", X_train_bow.get_shape()[1])

some feature names ['aa', 'aaa', 'aaaa', 'aaaaaaaaaaaa', 'aaaaaaaaaaaaaaaa',
'aaaaaaahhhhhh', 'aaaaaaawwwwwwww', 'aaaah', 'aaaand', 'aaah']
=====
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer (61441, 46268)
the number of unique words 46268
```

### [4.2] TF-IDF

```
In [396]: tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
tf_idf_vect.fit(X_train)
print("some sample features(unique words in the corpus)",tf_idf_vect.get_feature_names())
print('='*50)

X_train_tfidf = tf_idf_vect.transform(X_train)
X_test_tfidf = 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 number of unique words including both unigrams and bigrams ", X_train_tfidf.get_feature_names())

some sample features(unique words in the corpus) ['aa', 'aafco', 'abdominal',
'ability', 'able', 'able add', 'able buy', 'able chew', 'able drink', 'able eat']
=====
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer (61441, 36214)
the number of unique words including both unigrams and bigrams 36214
```

## Applying Multinomial Naive Bayes

### [5.1] Applying Naive Bayes on BOW, SET 1

```
In [397]: from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import confusion_matrix, roc_auc_score
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import roc_curve, auc
```

```

In [398]: # Function to find the best optimal Alpha using MultinomialNB and GridSearchCV
def getOptimizedAlpha(X_train_vector, Y_train):
    alpha_vals = np.logspace(-4, 4, 20) # 10^-4 to 10^4 and no. of samples =20

    param_grid = dict(alpha=alpha_vals)
    clf = MultinomialNB()

    grid = GridSearchCV(clf, param_grid, cv=10, scoring='roc_auc') # 10 fold cross validation
    grid.fit(X_train_vector, Y_train)

    print("Best cross-validation score: {:.2f}".format(grid.best_score_)) # best score
    best_alpha = round(grid.best_params_['alpha'],3) # best alpha value after 10 fold cross validation
    print("Best parameters: ", best_alpha)

    train_auc= grid.cv_results_['mean_train_score']
    train_auc_std= grid.cv_results_['std_train_score']
    cv_auc = grid.cv_results_['mean_test_score']
    cv_auc_std= grid.cv_results_['std_test_score']

    plt.plot(alpha_vals, train_auc, label='Train AUC')
    # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
    # create a shaded area between [mean - std, mean + std]
    plt.gca().fill_between(alphas, train_auc - train_auc_std, train_auc + train_auc_std, label='Train AUC')

    plt.plot(alpha_vals, cv_auc, label='CV AUC')
    # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
    plt.gca().fill_between(alphas, cv_auc - cv_auc_std, cv_auc + cv_auc_std, label='CV AUC')

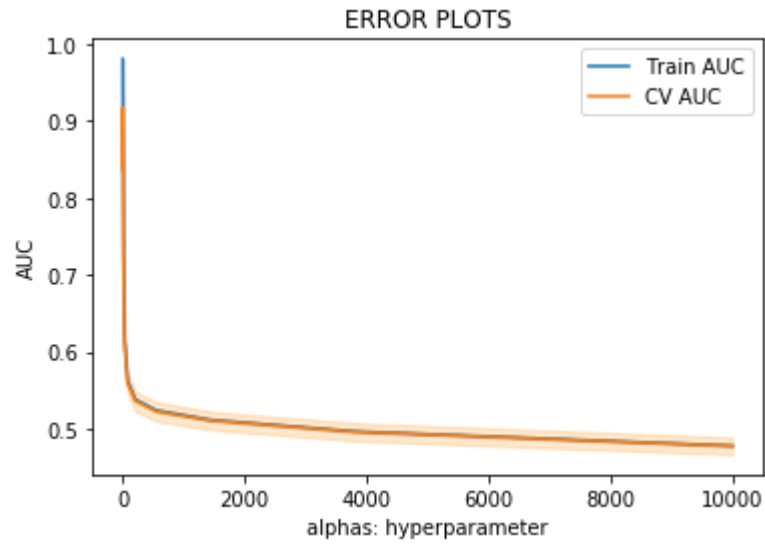
    plt.legend()
    plt.xlabel("alphas: hyperparameter")
    plt.ylabel("AUC")
    plt.title("ERROR PLOTS")
    plt.show()
    return best_alpha

```

```
In [399]: bestAlpha_bow = getOptimizedAlpha(X_train_bow, Y_train)
```

Best cross-validation score: 0.92

Best parameters: 0.234



```

In [400]: # Find the AUC for MultinomialNB
def NB_test(X_train_vector, Y_train, X_test_vector, Y_test, best_Alpha):
    clf = MultinomialNB(alpha=best_Alpha)
    clf.fit(X_train_vector, Y_train)

    train_fpr, train_tpr, threshold = roc_curve(Y_train, clf.predict_log_proba(X_train_vector))
    test_fpr, test_tpr, threshold = roc_curve(Y_test, clf.predict_log_proba(X_test_vector))

    test_auc = auc(test_fpr, test_tpr)
    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(test_auc))

    plt.legend()
    plt.xlabel("Alpha -> Hyperparameter")
    plt.ylabel("AUC")
    plt.title("Error Plot")
    plt.show()

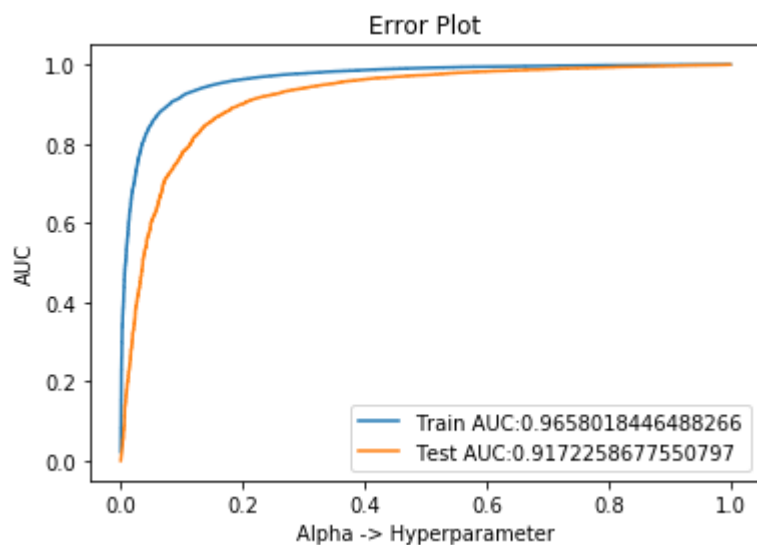
    print("training confusion matrix")
    print(confusion_matrix(Y_train, clf.predict(X_train_vector)))
    # print("Test confusion matrix")
    cm = confusion_matrix(Y_test, clf.predict(X_test_vector))
    # print(cm)

    # plot confusion matrix to describe the performance of classifier.
    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 of Test data")
    plt.xlabel("Predicted Label")
    plt.ylabel("True Label")
    plt.show()

    return test_auc

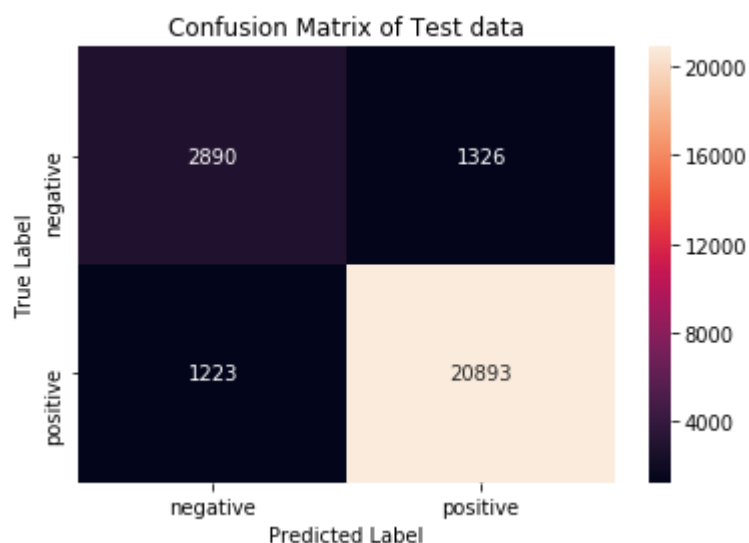
```

```
In [401]: test_auc_bow = NB_test(X_train_bow, Y_train, X_test_bow, Y_test, bestAlpha_bow)
```



training confusion matrix

```
[[ 8123 1842]  
 [ 2088 49388]]
```



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



```
In [402]: # Below is the vocabulary dataframe of each word in Bag of words vector with its
clf = MultinomialNB(alpha=bestAlpha_bow)
clf.fit(X_train_bow , Y_train)
features = count_vect.get_feature_names()
df = pd.DataFrame(clf.feature_log_prob_, columns=features)
df_new = df.T
print(df_new.head(3))
print('')
print('Feature Importance for the BOW')

# selecting the top 10 positive features using log_probabilities:
pos_features = df_new[1].sort_values(ascending = False)[0:10]
print(pos_features)
```

	0	1
aa	-12.798859	-11.848270
aaa	-14.461554	-13.060756
aaaa	-14.461554	-13.330184

Feature Importance for the BOW

not	-3.720183
like	-4.523675
good	-4.666490
great	-4.739032
one	-4.881107
taste	-4.962445
coffee	-5.014173
would	-5.065294
flavor	-5.074650
love	-5.076016

Name: 1, dtype: float64

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

```
In [403]: # selecting the top 10 negative features using log_probabilities:
neg_features = df_new[0].sort_values(ascending = False)[0:10]
print(neg_features)
```

not	-3.287500
like	-4.395847
would	-4.657690
taste	-4.690077
product	-4.709279
one	-4.881942
good	-5.133152
coffee	-5.143842
no	-5.165572
flavor	-5.209681

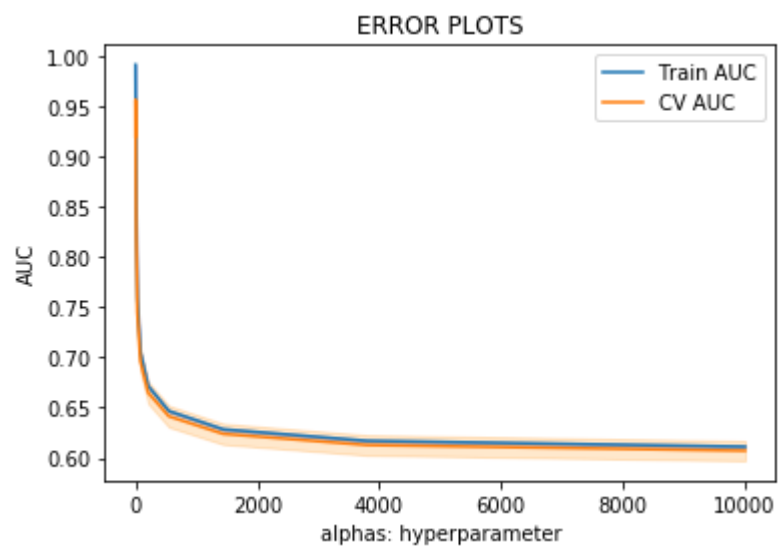
Name: 0, dtype: float64

## [5.2] Applying Naive Bayes on TFIDF, SET 2

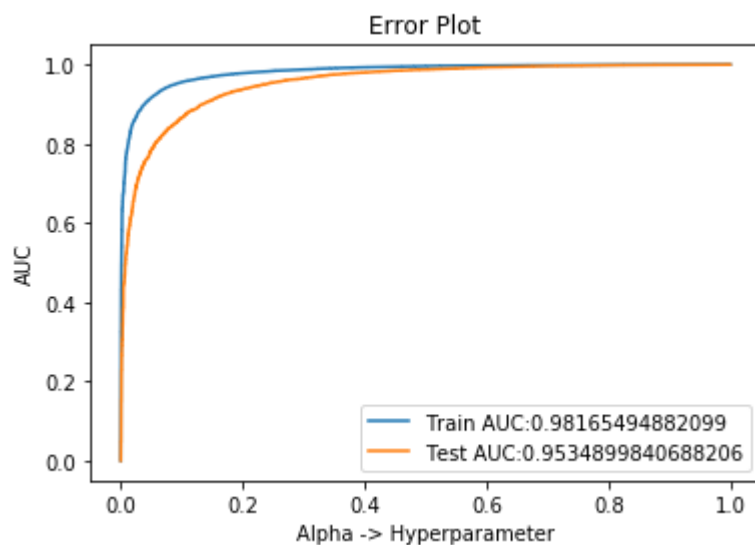
```
In [404]: bestAlpha_tfidf = getOptimizedAlpha(X_train_tfidf, Y_train)
```

Best cross-validation score: 0.96

Best parameters: 0.234

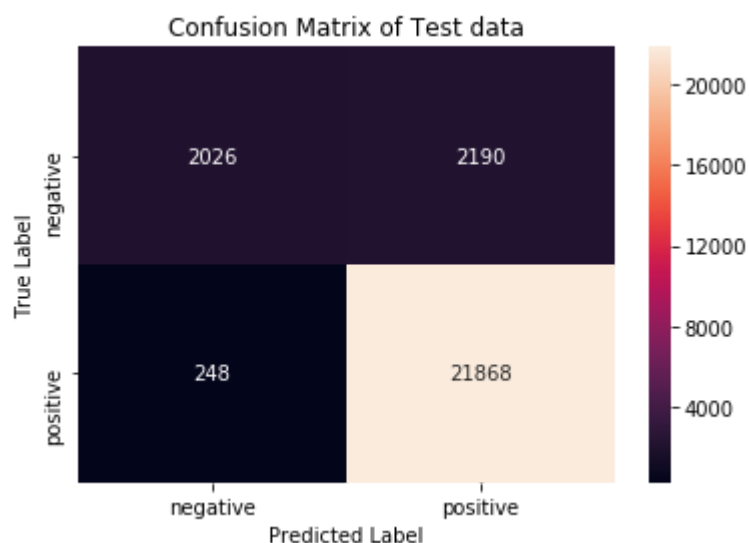


```
In [405]: test_auc_tfidf = NB_test(X_train_tfidf, Y_train, X_test_tfidf, Y_test, bestAlpha_
```



training confusion matrix

```
[[ 6285  3680]
 [  413 51063]]
```



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

```
In [406]: # Below is the vocabulary dataframe of each word in TFIDF vector with its log prob
clf = MultinomialNB(alpha=bestAlpha_tfidf)
clf.fit(X_train_tfidf , Y_train)
features = tf_idf_vect.get_feature_names()
df = pd.DataFrame(clf.feature_log_prob_, columns=features)
df_new = df.T
print(df_new.head(3))
print('')
print('Feature Importance for the TFIDF')

# selecting the top 10 positive features using log_probabilities:
pos_features = df_new[1].sort_values(ascending = False)[0:10]
print(pos_features)
```

	0	1
aa	-11.739227	-11.590272
aafco	-12.207235	-12.196399
abdominal	-11.807919	-12.271227

Feature Importance for the TFIDF

not	-5.319333
great	-5.658670
good	-5.727540
like	-5.777388
coffee	-5.826651
love	-5.903160
tea	-5.910888
one	-6.015894
taste	-6.025175
product	-6.031034

Name: 1, dtype: float64

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

```
In [407]: neg_features = df_new[0].sort_values(ascending = False)[0:10]
print(neg_features)
```

not	-4.841025
like	-5.649777
product	-5.761042
taste	-5.785964
would	-5.786284
coffee	-6.023545
one	-6.077367
no	-6.218726
flavor	-6.240156
good	-6.326414

Name: 0, dtype: float64

## [6] Conclusions

```
In [408]: models = pd.DataFrame({
'Vectorizer': ["BOW", "TFIDF"],
'Model' : ['MultinomialNB', 'MultinomialNB'],
'Hyper Parameter(Alpha)': [bestAlpha_bow, bestAlpha_tfidf],
'AUC': [test_auc_bow, test_auc_tfidf]},
columns = ["Vectorizer", "Model", "Hyper Parameter(Alpha)", "AUC"])
print(models)
```

	Vectorizer	Model	Hyper Parameter(Alpha)	AUC
0	BOW	MultinomialNB	0.234	0.917226
1	TFIDF	MultinomialNB	0.234	0.953490

- Hyper parameter 'Alpha' is calculated same i.e 0.234 by both BOW and TFIDF
- AUC of TFIDF is better than BOW
- Naive Bayes is super fast compared to KNN