

# 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 [5]: 1 %matplotlib inline
2 import warnings
3 warnings.filterwarnings("ignore")
4
5
6 import sqlite3
7 import pandas as pd
8 import numpy as np
9 import nltk
10 import string
11 import matplotlib.pyplot as plt
12 import seaborn as sns
13 from sklearn.feature_extraction.text import TfidfTransformer
14 from sklearn.feature_extraction.text import TfidfVectorizer
15
16 from sklearn.feature_extraction.text import CountVectorizer
17 from sklearn.metrics import confusion_matrix
18 from sklearn import metrics
19 from sklearn.metrics import roc_curve, auc
20 from nltk.stem.porter import PorterStemmer
21
22 import re
23 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
24 import string
25 from nltk.corpus import stopwords
26 from nltk.stem import PorterStemmer
27 from nltk.stem.wordnet import WordNetLemmatizer
28
29 from gensim.models import Word2Vec
30 from gensim.models import KeyedVectors
31 import pickle
32
33 from tqdm import tqdm
34 import os
```

```
C:\Users\sujpanda\Anaconda3\lib\site-packages\gensim\utils.py:1212: UserWarning:
g: detected Windows; aliasing chunkize to chunkize_serial
  warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")
```

```

In [101]: 1 # using SQLite Table to read data.
2 con = sqlite3.connect('C:\\Users\\sujpanda\\Desktop\\applied\\database.sqlite')
3
4 # filtering only positive and negative reviews i.e.
5 # not taking into consideration those reviews with Score=3
6 # SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000, will give top 500000 data points
7 # you can change the number to any other number based on your computing power
8
9 # filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3
10 # for tsne assignment you can take 5k data points
11
12 filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3
13
14 # Give reviews with Score>3 a positive rating(1), and reviews with a score<3
15 def partition(x):
16     if x < 3:
17         return 0
18     return 1
19
20 #changing reviews with score less than 3 to be positive and vice-versa
21 actualScore = filtered_data['Score']
22 positiveNegative = actualScore.map(partition)
23 filtered_data['Score'] = positiveNegative
24 print("Number of data points in our data", filtered_data.shape)
25 filtered_data.head(3)

```

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

Out[101]:

		Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian		1	
1	2	B00813GRG4	A1D87F6ZCVE5NK		dll pa	0	
2	3	B000LQOCH0	ABXLMWJIXXAIN		Natalia Corres "Natalia Corres"	1	

```

In [102]: 1 display = pd.read_sql_query("""
2 SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
3 FROM Reviews
4 GROUP BY UserId
5 HAVING COUNT(*)>1
6 """, con)

```

```
In [103]: 1 print(display.shape)
          2 display.head()
```

(80668, 7)

Out[103]:

	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
4	#oc-R12KPBODL2B5ZD	B007OSBE1U	Christopher P. Presta	1348617600	1	I didnt like this coffee. Instead of telling y...	2

```
In [104]: 1 display[display['UserId'] == 'AZY10LLTJ71NX']
```

Out[104]:

	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 [105]: 1 display['COUNT(*)'].sum()
```

Out[105]: 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 [106]: 1 display= pd.read_sql_query("""
2 SELECT *
3 FROM Reviews
4 WHERE Score != 3 AND UserId="AR5J8UI46CURR"
5 ORDER BY ProductID
6 """, con)
7 display.head()
```

Out[106]:

	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 [107]: 1 #Sorting data according to ProductId in ascending order
          2 sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True, in
```

```
In [108]: 1 #Deduplication of entries
          2 final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"
          3 final.shape
```

Out[108]: (87775, 10)

```
In [109]: 1 #Checking to see how much % of data still remains
          2 (final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

Out[109]: 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 [110]: 1 display= pd.read_sql_query("""
          2 SELECT *
          3 FROM Reviews
          4 WHERE Score != 3 AND Id=44737 OR Id=64422
          5 ORDER BY ProductID
          6 """, con)
          7
          8 display.head()
```

Out[110]:

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

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

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

```
In [111]: 1 final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [112]: 1 #Before starting the next phase of preprocessing Lets see the number of entri
          2 print(final.shape)
          3
          4 #How many positive and negative reviews are present in our dataset?
          5 final['Score'].value_counts()
```

```
(87773, 10)
```

```
Out[112]: 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 [113]: 1 # printing some random reviews
2 sent_0 = final['Text'].values[0]
3 print(sent_0)
4 print("="*50)
5
6 sent_1000 = final['Text'].values[1000]
7 print(sent_1000)
8 print("="*50)
9
10 sent_1500 = final['Text'].values[1500]
11 print(sent_1500)
12 print("="*50)
13
14 sent_4900 = final['Text'].values[4900]
15 print(sent_4900)
16 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.

=====

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 these 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 [114]: 1 # remove urls from text python: https://stackoverflow.com/a/40823105/4084039
2 sent_0 = re.sub(r"http\S+", "", sent_0)
3 sent_1000 = re.sub(r"http\S+", "", sent_1000)
4 sent_150 = re.sub(r"http\S+", "", sent_1500)
5 sent_4900 = re.sub(r"http\S+", "", sent_4900)
6
7 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 [115]: 1 # https://stackoverflow.com/questions/16206380/python-beautifulsoup-how-to-re
2 from bs4 import BeautifulSoup
3
4 soup = BeautifulSoup(sent_0, 'lxml')
5 text = soup.get_text()
6 print(text)
7 print("="*50)
8
9 soup = BeautifulSoup(sent_1000, 'lxml')
10 text = soup.get_text()
11 print(text)
12 print("="*50)
13
14 soup = BeautifulSoup(sent_1500, 'lxml')
15 text = soup.get_text()
16 print(text)
17 print("="*50)
18
19 soup = BeautifulSoup(sent_4900, 'lxml')
20 text = soup.get_text()
21 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 these 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 [116]: 1 # https://stackoverflow.com/a/47091490/4084039
2 import re
3
4 def decontracted(phrase):
5     # specific
6     phrase = re.sub(r"won't", "will not", phrase)
7     phrase = re.sub(r"can't", "can not", phrase)
8
9     # general
10    phrase = re.sub(r"n't", " not", phrase)
11    phrase = re.sub(r"\ 're", " are", phrase)
12    phrase = re.sub(r"\ 's", " is", phrase)
13    phrase = re.sub(r"\ 'd", " would", phrase)
14    phrase = re.sub(r"\ 'll", " will", phrase)
15    phrase = re.sub(r"\ 't", " not", phrase)
16    phrase = re.sub(r"\ 've", " have", phrase)
17    phrase = re.sub(r"\ 'm", " am", phrase)
18    return phrase
```

```
In [117]: 1 sent_1500 = decontracted(sent_1500)
2 print(sent_1500)
3 print("="*50)
```

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

```
In [118]: 1 #remove words with numbers python: https://stackoverflow.com/a/18082370/4084039
2 sent_0 = re.sub("\S*\d\S*", "", sent_0).strip()
3 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 [119]: 1 #remove spacial character: https://stackoverflow.com/a/5843547/4084039
2 sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)
3 print(sent_1500)
```

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

```
In [120]: 1 # https://gist.github.com/sebleier/554280
2 # we are removing the words from the stop words list: 'no', 'nor', 'not'
3 # <br /><br /> ==> after the above steps, we are getting "br br"
4 # we are including them into stop words list
5 # instead of <br /> if we have <br/> these tags would have revmoved in the 1s
6
7 stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours',
8               "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he',
9               'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'i',
10              'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'is',
11              'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have',
12              'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'beca',
13              'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into',
14              'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on',
15              'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 't',
16              'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'th',
17              's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "shoul",
18              've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn',
19              "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'm',
20              "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shou",
21              'won', "won't", 'wouldn', "wouldn't"])
```

```
In [121]: 1 # Combining all the above stundents
2 from tqdm import tqdm
3 preprocessed_reviews = []
4 # tqdm is for printing the status bar
5 for sentence in tqdm(final['Text'].values):
6     sentence = re.sub(r"http\S+", "", sentence)
7     sentence = BeautifulSoup(sentence, 'lxml').get_text()
8     sentence = decontracted(sentence)
9     sentence = re.sub("\S*\d\S*", "", sentence).strip()
10    sentence = re.sub('[^A-Za-z]+', ' ', sentence)
11    # https://gist.github.com/sebleier/554280
12    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not
13    preprocessed_reviews.append(sentence.strip())
```

100%|██████████| 87773/87773 [01:18<00:00, 1111.85it/s]

```
In [122]: 1 preprocessed_reviews[1500]
```

Out[122]: 'way hot blood took bite jig lol'

## [3.2] Preprocessing Review Summary

```
In [123]: 1 from tqdm import tqdm
2 preprocessed_summary = []
3 # tqdm is for printing the status bar
4 for sentence in tqdm(final['Summary'].values):
5     sentence = re.sub(r"http\S+", "", sentence)
6     sentence = BeautifulSoup(sentence, 'lxml').get_text()
7     sentence = decontracted(sentence)
8     sentence = re.sub("\S*\d\S*", "", sentence).strip()
9     sentence = re.sub('[^A-Za-z]+', ' ', sentence)
10    # https://gist.github.com/sebleier/554280
11    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not
12    preprocessed_summary.append(sentence.strip())
```

100%|██████████| 87773/87773 [01:01<00:00, 1434.16it/s]

```
In [124]: 1 print(preprocessed_summary[1000])
```

not much taste

## [5] Assignment 4: Apply Naive Bayes

### 1. Apply Multinomial NaiveBayes on these feature sets

- **SET 1:** Review text, preprocessed one converted into vectors using (BOW)
- **SET 2:** Review text, preprocessed one converted into vectors using (TFIDF)

### 2. The hyper paramter tuning(find best Alpha)

- Find the best hyper parameter which will give the maximum [AUC](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/>) value
- Consider a wide range of alpha values for hyperparameter tuning, start as low as 0.00001
- 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. Feature importance


- Find the top 10 features of positive class and top 10 features of negative class for both feature sets **Set 1** and **Set 2** using absolute values of `coef_` parameter of [MultinomialNB](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html) ([https://scikit-learn.org/stable/modules/generated/sklearn.naive\\_bayes.MultinomialNB.html](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html)) and print their corresponding feature names


### 4. 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.

## 5. 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. Here on X-axis you will have alpha values, since they have a wide range, just to represent those alpha values on the graph, apply log function on those alpha values.

 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](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/) (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/>) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](https://seaborn.pydata.org/generated/seaborn.heatmap.html).

 (<https://seaborn.pydata.org/generated/seaborn.heatmap.html>)

## 6. 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](http://zetcode.com/python/prettytable/) (<http://zetcode.com/python/prettytable/>)



### Note: Data Leakage

- There will be an issue of data-leakage if you vectorize the entire data and then split it into train/cv/test.
- To avoid the issue of data-leakag, make sure to split your data first and then vectorize it.
- While vectorizing your data, apply the method `fit_transform()` on you train data, and apply the method `transform()` on cv/test data.
- For more details please go through this [link](https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf). (<https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf>)

In [125]:

```

1  ## Some utility functions
2
3  def get_optimum_alpha(X_train,X_test,y_train,y_test):
4
5      from sklearn.metrics import roc_curve
6      from sklearn.metrics import roc_auc_score
7
8      alphas_range1 = ['0.000001','0.00001','0.0001','0.001','0.01','0.1','10',
9      alphas_range = [0.000001,0.00001,0.0001,0.00,0.01,0.1,10,100,1000,10000,1
10     dummy_range = [1,2,3,4,5,6,7,8,9,10,11]
11
12     auc_scores = []
13     auc_train_scores = []
14
15     i = 0
16     for i in alphas_range:
17         clf = MultinomialNB(alpha=i)
18
19         # fitting the model on crossvalidation train
20         clf.fit(X_train, y_train)
21
22
23         #evaluate AUC score.
24         probs = clf.predict_proba(X_test)
25         probs = probs[:, 1]
26         # calculate AUC
27         auc = roc_auc_score(y_test, probs)
28         print('AUC: %.3f' % auc)
29         auc_scores.append(auc)
30
31     print('#####')
32     print('AUC from train data #####')
33     i = 0
34     for i in alphas_range:
35         clf = MultinomialNB(i)
36
37         # fitting the model on crossvalidation train
38         clf.fit(X_train, y_train)
39
40         #evaluate AUC score.
41         probs = clf.predict_proba(X_train)
42         probs = probs[:, 1]
43         # calculate AUC
44         auc = roc_auc_score(y_train, probs)
45         print('AUC: %.3f' % auc)
46         auc_train_scores.append(auc)
47
48     plt.plot(dummy_range, auc_scores,'r')
49     plt.plot(dummy_range, auc_train_scores,'b')
50     plt.xticks(dummy_range, alphas_range1, rotation='vertical')
51     for xy in zip(dummy_range, auc_scores):
52         plt.annotate('(%f, %f)' % xy, xy=xy, textcoords='data')
53     for xy in zip(dummy_range, auc_train_scores):
54         plt.annotate('(%f, %f)' % xy, xy=xy, textcoords='data')
55
56

```

```

57 plt.xlabel('Alphas')
58 plt.ylabel('auc_scores')
59 plt.show()
60

```

```

In [126]: 1 def nb_results(optimum_alpha,X_train,X_test,y_train,y_test):
2         # roc curve and auc
3         from sklearn.datasets import make_classification
4         from sklearn.model_selection import train_test_split
5         from sklearn.metrics import roc_curve
6         from sklearn.metrics import roc_auc_score
7         from matplotlib import pyplot
8         # ===== KNN with k = optimal_k =====
9         # instantiate learning model with optimum alpha
10        clf = MultinomialNB(alpha=optimum_alpha)
11
12        # fitting the model
13        clf.fit(X_train, y_train)
14
15        # predict the response
16        pred = clf.predict(X_test)
17
18        # evaluate accuracy
19        acc = accuracy_score(y_test, pred) * 100
20        print('\nThe accuracy of the NB classifier for alpha = %f is %f%%' % (opt
21
22        probs = clf.predict_proba(X_test)
23        probs = probs[:, 1]
24        # calculate AUC
25        auc = roc_auc_score(y_test, probs)
26        print('AUC: %.3f' % auc)
27        # calculate roc curve
28        fpr, tpr, thresholds = roc_curve(y_test, probs)
29        # plot no skill
30        pyplot.plot([0, 1], [0, 1], linestyle='--')
31        # plot the roc curve for the model
32        pyplot.plot(fpr, tpr, marker='.')
33        # show the plot
34        pyplot.show()
35        from sklearn.metrics import confusion_matrix
36        con_mat = confusion_matrix(y_test, pred, [0, 1])
37        return con_mat

```

```

In [127]: 1 def showHeatMap(con_mat):
2         class_label = ["negative", "positive"]
3         df_cm = pd.DataFrame(con_mat, index = class_label, columns = class_label)
4         sns.heatmap(df_cm, annot = True, fmt = "d")
5         plt.title("Confusion Matrix")
6         plt.xlabel("Predicted Label")
7         plt.ylabel("True Label")
8         plt.show()

```

```
In [128]: 1 def most_informative_feature_for_class(vectorizer, classifier, classlabel, n=  
2         labelid = list(classifier.classes_).index(classlabel)  
3         feature_names = vectorizer.get_feature_names()  
4         topn = sorted(zip(classifier.coef_[labelid], feature_names))[-n:]  
5  
6         for coef, feat in topn:  
7             print(classlabel, feat, coef)
```

## Applying Multinomial Naive Bayes

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

```
In [129]: 1 from sklearn.cross_validation import train_test_split  
2 from sklearn.naive_bayes import MultinomialNB  
3 from sklearn.metrics import accuracy_score  
4 from sklearn.cross_validation import cross_val_score  
5 from collections import Counter  
6 from sklearn.metrics import accuracy_score  
7 from sklearn import cross_validation  
8 import warnings  
9 warnings.filterwarnings("ignore")
```

```
In [130]: 1 print(final['Text'].shape)  
  
(87773,)
```

```
In [131]: 1 X_1, X_test, y_1, y_test = cross_validation.train_test_split(preprocessed_rev  
2  
◀────────────────────────────────────────────────────────────────────────────────▶
```



```

In [132]: 1 count_vect = CountVectorizer()
          2 final_counts = count_vect.fit_transform(X_1)
          3 final_test_count = count_vect.transform(X_test)
          4
          5 # split the train data set into cross validation train and cross validation t
          6 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_1, y_1, test_siz
          7
          8 final_counts_tr_cv = count_vect.transform(X_tr)
          9 final_test_count_cv = count_vect.transform(X_cv)
         10
         11 get_optimum_alpha(final_counts_tr_cv,final_test_count_cv,y_tr,y_cv)

```

AUC: 0.783

AUC: 0.802

AUC: 0.826

AUC: 0.741

AUC: 0.884

AUC: 0.907

AUC: 0.697

AUC: 0.548

AUC: 0.517

AUC: 0.477

AUC: 0.443

#####  
 AUC from train data #####

AUC: 0.984

AUC: 0.983

AUC: 0.983

AUC: 0.984

AUC: 0.979

AUC: 0.972

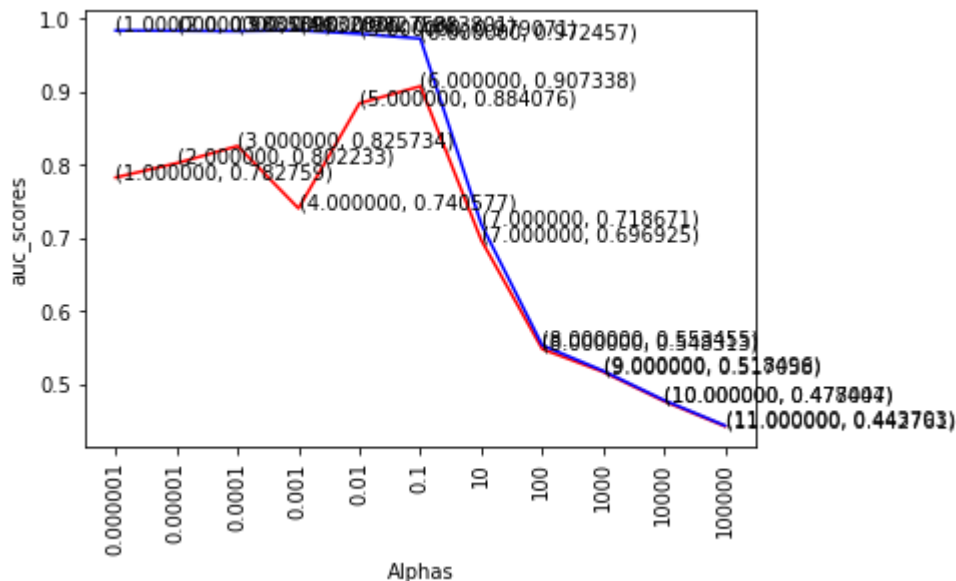
AUC: 0.719

AUC: 0.553

AUC: 0.518

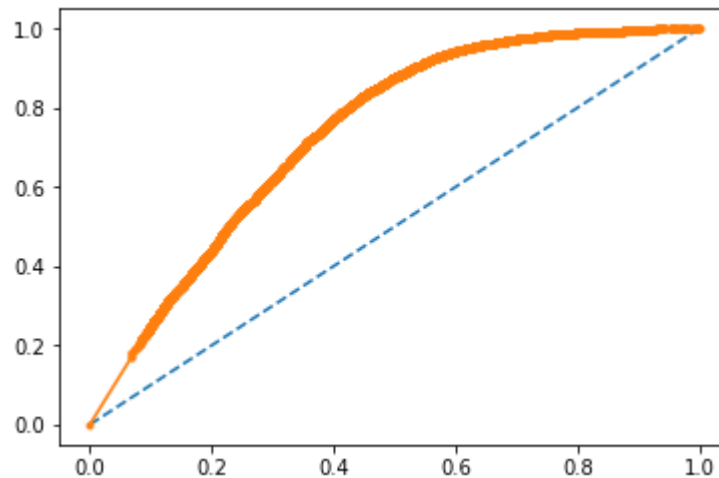
AUC: 0.478

AUC: 0.444



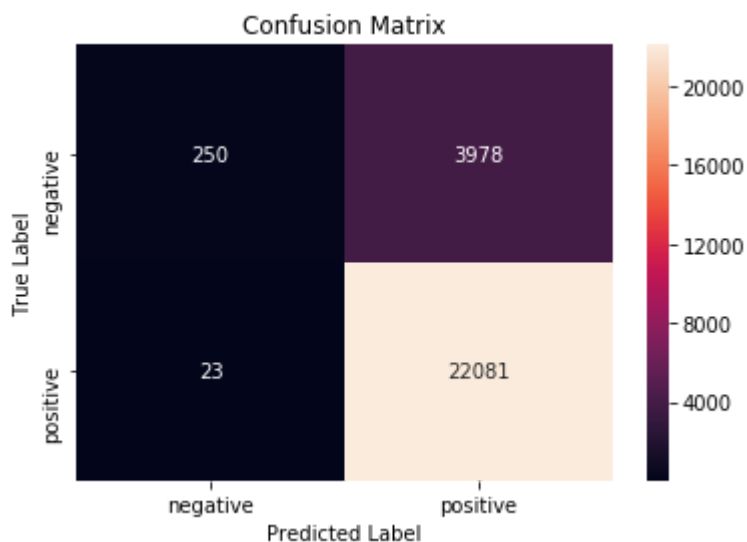
```
In [133]: 1 con_mat=nb_results(10,final_counts,final_test_count,y_1,y_test)
```

The accuracy of the NB classifier for alpha = 10.000000 is 84.805560%  
AUC: 0.738



Observation : With alpha = 10 the accuracy is 86 % and AUC is 0.738 which is much better than dumb model.

```
In [134]: 1 showHeatMap(con_mat)
```



Observation : My model misclassified 23 + 3978 points

```
In [135]: 1 bow_featurennames = count_vect.get_feature_names()
          2 print(bow_featurennames[100])
```

abruzzo

```
In [136]: 1 clf = MultinomialNB(alpha=10)
          2 clf.fit(final_counts, y_1)
          3 pred = clf.predict(final_test_count)
          4
```

```
In [137]: 1 feat_count = clf.feature_count_
          2 feat_count.shape
```

Out[137]: (2, 46446)

```
In [138]: 1 # Number of samples encountered for each class during fitting
          2 clf.class_count_
```

Out[138]: array([ 9953., 51488.])

```
In [139]: 1 probs = clf.feature_log_prob_
```

```
In [140]: 1 feature_prob = pd.DataFrame(probs, columns = bow_featurenames)
          2 feature_prob_tr = feature_prob.T
          3 feature_prob_tr.shape
```

Out[140]: (46446, 2)

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

```
In [141]: 1 print("\n\n Top 10 Positive Features:-\n", feature_prob_tr[1].sort_values(asc=
```

```
Top 10 Positive Features:-
not      -3.925894
like     -4.733360
good     -4.864065
great    -4.951355
one      -5.084318
taste    -5.172942
coffee  -5.206030
flavor   -5.275124
love     -5.281123
would    -5.282327
Name: 1, dtype: float64
```

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

```
In [142]: 1 print("Top 10 Negative Features:-\n",feature_prob_tr[0].sort_values(ascending
Top 10 Negative Features:-
not      -3.999070
like     -5.118756
would    -5.353858
product  -5.389252
taste    -5.400513
one       -5.598615
coffee   -5.823486
good      -5.832546
flavor    -5.882702
no        -5.886297
Name: 0, dtype: float64
```

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

```
In [143]: 1 X_1, X_test, y_1, y_test = cross_validation.train_test_split(preprocessed_rev
```

In [144]:

```

1 tf_idf_vect = TfidfVectorizer()
2 tf_idf_vect.fit(X_1)
3 final_tf_idf = tf_idf_vect.transform(X_1)
4 final_test_count = tf_idf_vect.transform(X_test)
5
6 # split the train data set into cross validation train and cross validation t
7 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_1, y_1, test_siz
8
9 final_counts_tr_cv = tf_idf_vect.transform(X_tr)
10 final_test_count_cv = tf_idf_vect.transform(X_cv)
11
12 get_optimum_alpha(final_counts_tr_cv, final_test_count_cv, y_tr, y_cv)

```

AUC: 0.771

AUC: 0.791

AUC: 0.818

AUC: 0.731

AUC: 0.890

AUC: 0.908

AUC: 0.735

AUC: 0.673

AUC: 0.635

AUC: 0.617

AUC: 0.614

#####

AUC from train data #####

AUC: 0.989

AUC: 0.988

AUC: 0.988

AUC: 0.989

AUC: 0.985

AUC: 0.975

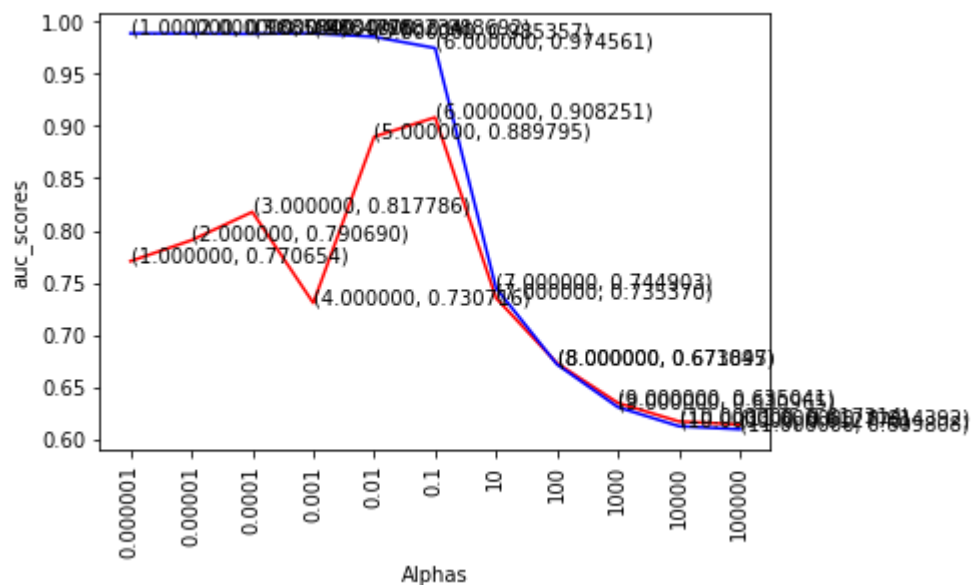
AUC: 0.745

AUC: 0.672

AUC: 0.631

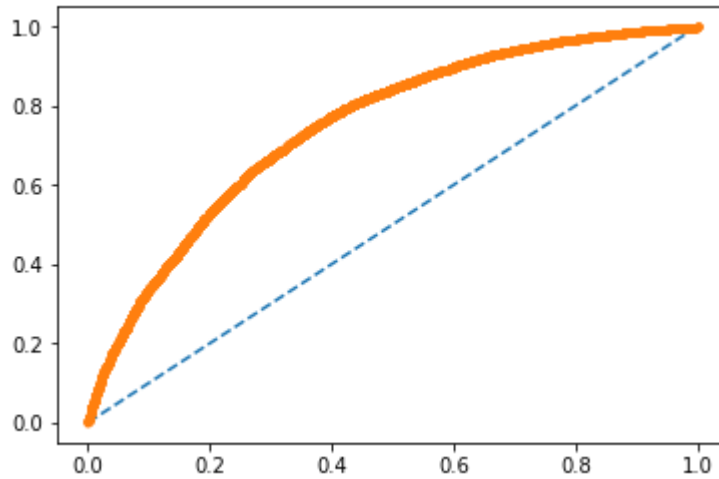
AUC: 0.613

AUC: 0.610



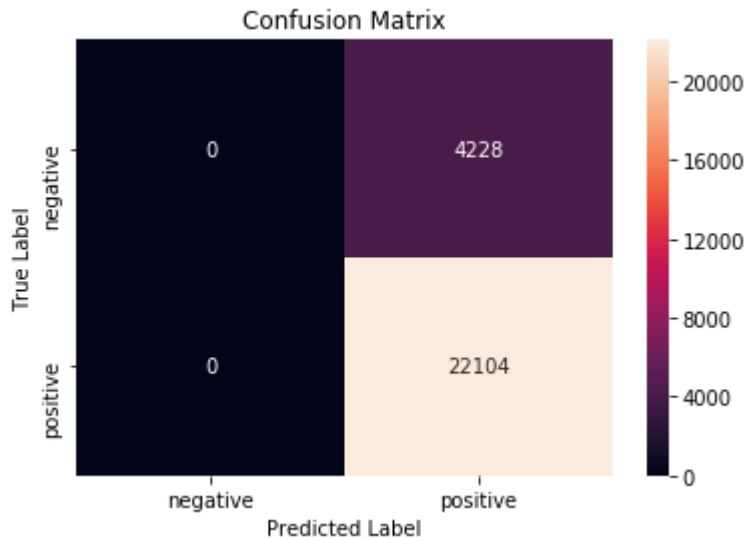
```
In [145]: 1 con_mat=nb_results(10,final_tf_idf,final_test_count,y_1,y_test)
```

The accuracy of the NB classifier for alpha = 10.000000 is 83.943491%  
AUC: 0.748



Observation : With alpha = 10 the accuracy is 83 % and AUC is 0.748.

```
In [146]: 1 showHeatMap(con_mat)
```



Observation: My model pedicted 4228 points wrongly

```
In [147]: 1 tfidf_featureennames = tf_idf_vect.get_feature_names()
```

```
In [148]: 1 clf = MultinomialNB(alpha=10)
2 clf.fit(final_tf_idf, y_1)
3 pred = clf.predict(final_test_count)
```

```
In [149]: 1 feat_count = clf.feature_count_
          2 feat_count.shape
```

```
Out[149]: (2, 46446)
```

```
In [150]: 1 clf.class_count_
```

```
Out[150]: array([ 9953., 51488.])
```

```
In [151]: 1 probs = clf.feature_log_prob_
```

```
In [152]: 1 feature_prob = pd.DataFrame(probs, columns = tfidf_featurenames)
          2 feature_prob_tr = feature_prob.T
          3 feature_prob_tr.shape
```

```
Out[152]: (46446, 2)
```

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

```
In [153]: 1 print("\n\n Top 10 Positive Features:-\n",feature_prob_tr[1].sort_values(ascending=False))
```

```
Top 10 Positive Features:-
not      -5.893515
great    -6.200917
good     -6.266220
coffee  -6.325771
like     -6.333298
love     -6.434347
tea      -6.456825
taste    -6.572348
one      -6.583339
product  -6.590492
Name: 1, dtype: float64
```

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

```
In [154]: 1 print("Top 10 Negative Features:-\n",feature_prob_tr[0].sort_values(ascending=False))
```

```
Top 10 Negative Features:-
not      -6.613617
like     -7.405604
taste    -7.503258
product  -7.516003
would    -7.533082
coffee  -7.710883
one      -7.837948
flavor   -7.917743
no       -7.993194
good     -8.041537
Name: 0, dtype: float64
```

In [ ]: 1

**repeat with extra features**

```
In [179]: 1 mylen = np.vectorize(len)
          2 newarr = mylen(preprocessed_summary)
```

```
In [180]: 1 newproce_reviews = np.asarray(preprocessed_reviews)
```

```
In [181]: 1 newproce_summary = np.asarray(preprocessed_summary)
```

```
In [182]: 1 df = pd.DataFrame({'desc':newproce_reviews, 'summary':newproce_summary, 'len':
```

```
In [183]: 1 df.head()
```

Out[183]:

	desc	summary	len
0	dogs loves chicken product china wont buying a...	made china	10
1	dogs love saw pet store tag attached regarding...	dog lover delites	17
2	infestation fruitflies literally everywhere fl...	one fruitfly stuck	18
3	worst product gotten long time would rate no s...	not work not waste money	24
4	wish would read reviews making purchase basica...	big rip	7

```
In [184]: 1 #df = df[:60000]
          2 print(len(final['Score']))
          3
```

87773

```
In [185]: 1 X_1, X_test, y_1, y_test = cross_validation.train_test_split(df, final['Score
```



In [186]:

```

1 import scipy
2 count_vect = CountVectorizer()
3 final_counts = count_vect.fit_transform(X_1['desc'])
4 final_test_count = count_vect.transform(X_test['desc'])
5
6 # split the train data set into cross validation train and cross validation test
7 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_1, y_1, test_size=0.2)
8
9 final_counts_tr_cv = count_vect.transform(X_tr['desc'])
10 final_test_count_cv = count_vect.transform(X_cv['desc'])
11
12 from scipy.sparse import csr_matrix, issparse
13
14 #####Adding len as feature#####
15 #if issparse(final_counts_tr_cv):
16     #print('sparse matrix')
17 len_sparse = scipy.sparse.coo_matrix(X_tr['len'])
18 len_sparse = len_sparse.transpose()
19
20 final_counts_tr_cv = scipy.sparse.hstack([final_counts_tr_cv, len_sparse])
21 print(final_counts_tr_cv.shape)
22
23 len_test_sparse = scipy.sparse.coo_matrix(X_cv['len'])
24 len_test_sparse = len_test_sparse.transpose()
25 final_test_count_cv = scipy.sparse.hstack([final_test_count_cv, len_test_sparse])
26 print("final_counts_tr_cv.shape after length = ", final_counts_tr_cv.shape)
27
28 #####Adding summary as feature#####
29 final_summary_count = count_vect.transform(X_tr['summary'])
30 final_test_summary_count_cv = count_vect.transform(X_cv['summary'])
31 columns=count_vect.get_feature_names()
32
33 print("sujet", final_summary_count[:,12].shape)
34
35
36 #df1 = pd.DataFrame(final_summary_count.toarray(), columns=count_vect.get_feature_names())
37 #print(df1['abandon'].shape)
38
39 #f1_sparse = scipy.sparse.coo_matrix(df1['abandon'])
40 #f1_sparse = f1_sparse.transpose()
41 final_counts_tr_cv = scipy.sparse.hstack([final_counts_tr_cv, final_summary_count])
42 print("final_counts_tr_cv.shape after f1= ", final_counts_tr_cv.shape)
43
44
45 #df2 = pd.DataFrame(final_test_summary_count_cv.toarray(), columns=count_vect.get_feature_names())
46
47 #f1_test_sparse = scipy.sparse.coo_matrix(df2['abandon'])
48 #f1_test_sparse = f1_test_sparse.transpose()
49 final_test_count_cv = scipy.sparse.hstack([final_test_count_cv, final_test_summary_count_cv])
50
51
52 #f1_sparse = scipy.sparse.coo_matrix(df1['zucchini'])
53 #f1_sparse = f1_sparse.transpose()
54 final_counts_tr_cv = scipy.sparse.hstack([final_counts_tr_cv, final_summary_count])
55 print("final_counts_tr_cv.shape after f2= ", final_counts_tr_cv.shape)
56

```

```

57
58 #df2 = pd.DataFrame(final_test_summary_count_cv.toarray(), columns=count_vect
59
60 #f1_test_sparse = scipy.sparse.coo_matrix(df2['zucchini'])
61 #f1_test_sparse = f1_test_sparse.transpose()
62 final_test_count_cv = scipy.sparse.hstack([final_test_count_cv, final_test_sum
63
64
65 get_optimum_alpha(final_counts_tr_cv, final_test_count_cv, y_tr, y_cv)

```

```

(43008, 46447)
final_counts_tr_cv.shape after length = (43008, 46447)
sujet (43008, 1)
final_counts_tr_cv.shape after f1= (43008, 46448)
final_counts_tr_cv.shape after f2= (43008, 46449)

```

AUC: 0.784

AUC: 0.804

AUC: 0.828

AUC: 0.740

AUC: 0.885

AUC: 0.906

AUC: 0.696

AUC: 0.529

AUC: 0.510

AUC: 0.502

AUC: 0.502

#####

AUC from train data #####

AUC: 0.983

AUC: 0.983

AUC: 0.983

AUC: 0.984

AUC: 0.979

AUC: 0.971

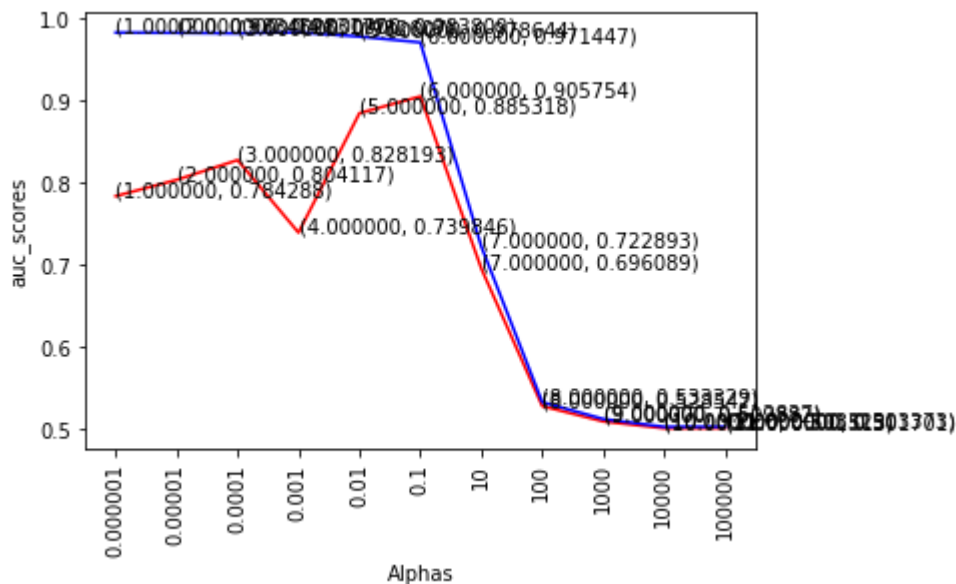
AUC: 0.723

AUC: 0.533

AUC: 0.513

AUC: 0.503

AUC: 0.503





```

In [187]: 1 print("Initial final countes = ",final_counts.shape)
2
3 #####Adding length as a feature#####
4 train_sparse = scipy.sparse.coo_matrix(X_1['len'])
5 train_sparse = train_sparse.transpose()
6 final_counts = scipy.sparse.hstack([final_counts,train_sparse])
7
8 test_sparse = scipy.sparse.coo_matrix(X_test['len'])
9 test_sparse = test_sparse.transpose()
10 final_test_count = scipy.sparse.hstack([final_test_count,test_sparse])
11
12 print("final_test_count.shape after length = ",final_test_count.shape)
13 print("Initial final countes after length = ",final_counts.shape)
14
15 #####Adding summary feature as a feature#####
16
17 final_summary = count_vect.transform(X_1['summary'])
18
19 #df3 = pd.DataFrame(final_summary.toarray(), columns=count_vect.get_feature_n
20
21 #print(df1['abandon'].shape)
22
23 #f1_sparse = scipy.sparse.coo_matrix(df3['abandon'])
24 #f1_sparse = f1_sparse.transpose()
25 final_counts = scipy.sparse.hstack([final_counts, final_summary[:,12]])
26
27 final_test_summary = count_vect.transform(X_test['summary'])
28 #df4 = pd.DataFrame(final_test_summary.toarray(), columns=count_vect.get_feat
29 #print(df1['abandon'].shape)
30
31 #f1_test_sparse = scipy.sparse.coo_matrix(df4['abandon'])
32 #f1_test_sparse = f1_test_sparse.transpose()
33 final_test_count = scipy.sparse.hstack([final_test_count, final_test_summary[
34 print("final_test_count.shape after f1 = ",final_test_count.shape)
35 print("Initial final countes after f1 = ",final_counts.shape)
36
37
38 #f1_sparse = scipy.sparse.coo_matrix(df3['zucchini'])
39 #f1_sparse = f1_sparse.transpose()
40 final_counts = scipy.sparse.hstack([final_counts, final_summary[:,112]])
41
42
43 #final_test_summary = count_vect.transform(X_test['summary'])
44 #df4 = pd.DataFrame(final_test_summary.toarray(), columns=count_vect.get_feat
45 #print(df1['abandon'].shape)
46
47 #f1_test_sparse = scipy.sparse.coo_matrix(df4['zucchini'])
48 #f1_test_sparse = f1_test_sparse.transpose()
49 final_test_count = scipy.sparse.hstack([final_test_count, final_test_summary[
50 print("final_test_count.shape after f2 = ",final_test_count.shape)
51 print("Initial final countes after f2 = ",final_counts.shape)
52
53
54
55 con_mat=nb_results(10,final_counts,final_test_count,y_1,y_test)

```

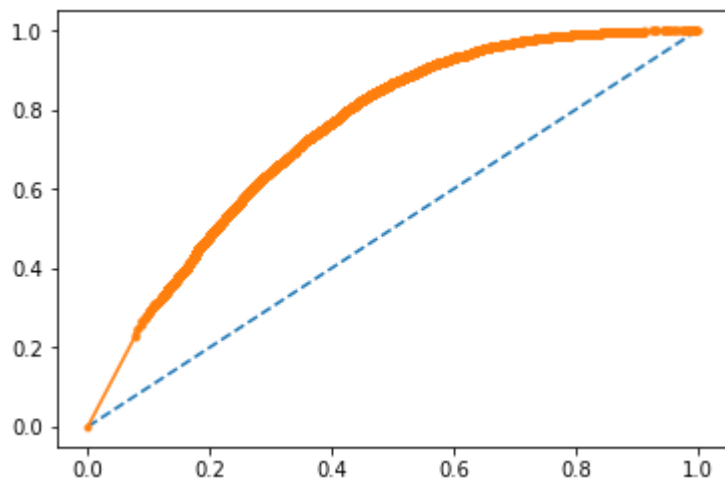
Initial final countes = (61441, 46446)

```

final_test_count.shape after length = (26332, 46447)
Initial final countes after length = (61441, 46447)
final_test_count.shape after f1 = (26332, 46448)
Initial final countes after f1 = (61441, 46448)
final_test_count.shape after f2 = (26332, 46449)
Initial final countes after f2 = (61441, 46449)

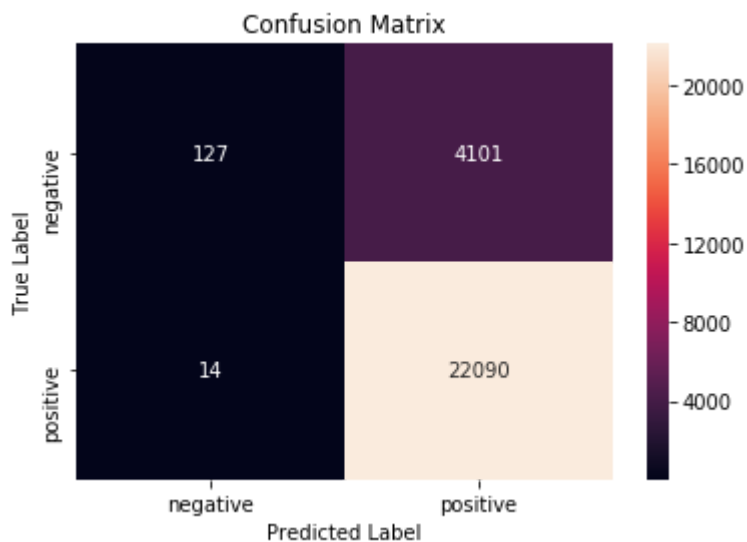
```

The accuracy of the NB classifier for alpha = 10.000000 is 84.372626%  
 AUC: 0.745



Observation : With alpha = 10 the accuracy is 84 % and AUC is 0.745 which is much better than dumb model.

In [188]: 1 showHeatMap(con\_mat)



Observation: My model pedicted 14 + 4101 points wrongly

In [198]: 1 X\_1, X\_test, y\_1, y\_test = cross\_validation.train\_test\_split(df, final['Score

In [199]:

```

1 tf_idf_vect = TfidfVectorizer()
2 tf_idf_vect.fit(X_1['desc'])
3 final_tf_idf = tf_idf_vect.transform(X_1['desc'])
4 final_test_count = tf_idf_vect.transform(X_test['desc'])
5
6 # split the train data set into cross validation train and cross validation test
7 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_1, y_1, test_size=0.2)
8
9 final_counts_tr_cv = tf_idf_vect.transform(X_tr['desc'])
10 final_test_count_cv = tf_idf_vect.transform(X_cv['desc'])
11
12 #####Adding len as feature#####
13 #if issparse(final_counts_tr_cv):
14     #print('sparse matrix')
15 len_sparse = scipy.sparse.coo_matrix(X_tr['len'])
16 len_sparse = len_sparse.transpose()
17
18 final_counts_tr_cv = scipy.sparse.hstack([final_counts_tr_cv, len_sparse])
19
20 len_test_sparse = scipy.sparse.coo_matrix(X_cv['len'])
21 len_test_sparse = len_test_sparse.transpose()
22 final_test_count_cv = scipy.sparse.hstack([final_test_count_cv, len_test_sparse])
23 print("final_counts_tr_cv.shape after len= ", final_counts_tr_cv.shape)
24
25 #####Adding summary as feature#####
26 final_summary_count = tf_idf_vect.transform(X_tr['summary'])
27 final_test_summary_count_cv = tf_idf_vect.transform(X_cv['summary'])
28 columns=tf_idf_vect.get_feature_names()
29
30
31 #df1 = pd.DataFrame(final_summary_count.toarray(), columns=columns)
32 #print(df1['abandon'].shape)
33
34 #f1_sparse = scipy.sparse.coo_matrix(df1['abandon'])
35 #f1_sparse = f1_sparse.transpose()
36 final_counts_tr_cv = scipy.sparse.hstack([final_counts_tr_cv, final_summary_count])
37 print("final_counts_tr_cv.shape after f1= ", final_counts_tr_cv.shape)
38
39
40 #df2 = pd.DataFrame(final_test_summary_count_cv.toarray(), columns=columns)
41
42 #f1_test_sparse = scipy.sparse.coo_matrix(df2['abandon'])
43 #f1_test_sparse = f1_test_sparse.transpose()
44 final_test_count_cv = scipy.sparse.hstack([final_test_count_cv, final_test_summary_count_cv])
45
46
47 #f1_sparse = scipy.sparse.coo_matrix(df1['zucchini'])
48 #f1_sparse = f1_sparse.transpose()
49 final_counts_tr_cv = scipy.sparse.hstack([final_counts_tr_cv, final_summary_count])
50 print("final_counts_tr_cv.shape after f2= ", final_counts_tr_cv.shape)
51
52
53 #df2 = pd.DataFrame(final_test_summary_count_cv.toarray(), columns=columns)
54
55 #f1_test_sparse = scipy.sparse.coo_matrix(df2['zucchini'])
56 #f1_test_sparse = f1_test_sparse.transpose()

```

```

57 final_test_count_cv = scipy.sparse.hstack([final_test_count_cv,final_test_sum
58
59 get_optimum_alpha(final_counts_tr_cv,final_test_count_cv,y_tr,y_cv)

```

final\_counts\_tr\_cv.shape after len= (43008, 46447)

final\_counts\_tr\_cv.shape after f1= (43008, 46448)

final\_counts\_tr\_cv.shape after f2= (43008, 46449)

AUC: 0.762

AUC: 0.783

AUC: 0.811

AUC: 0.720

AUC: 0.885

AUC: 0.903

AUC: 0.548

AUC: 0.519

AUC: 0.514

AUC: 0.513

AUC: 0.512

#####

AUC from train data #####

AUC: 0.989

AUC: 0.989

AUC: 0.988

AUC: 0.989

AUC: 0.985

AUC: 0.970

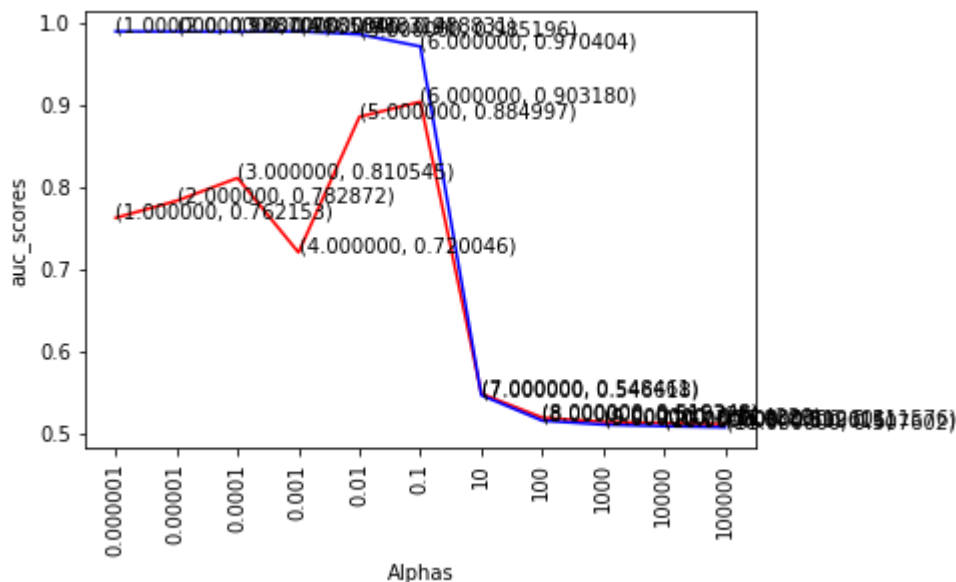
AUC: 0.547

AUC: 0.516

AUC: 0.511

AUC: 0.509

AUC: 0.508



```

In [200]: 1 train_sparse = scipy.sparse.coo_matrix(X_1['len'])
2 train_sparse = train_sparse.transpose()
3 final_tf_idf = scipy.sparse.hstack([final_tf_idf,train_sparse])
4
5 test_sparse = scipy.sparse.coo_matrix(X_test['len'])
6 test_sparse = test_sparse.transpose()
7 final_test_count = scipy.sparse.hstack([final_test_count,test_sparse])
8
9 #####Adding summary feature as a feature#####
10 final_summary = count_vect.transform(X_1['summary'])
11
12 #df3 = pd.DataFrame(final_summary.toarray(), columns=tf_idf_vect.get_feature_
13
14 #f1_sparse = scipy.sparse.coo_matrix(df3['abandon'])
15 #f1_sparse = f1_sparse.transpose()
16 final_tf_idf = scipy.sparse.hstack([final_tf_idf, final_summary[:,12]])
17 print("final_counts_tr_cv.shape = ",final_tf_idf.shape)
18
19 final_test_summary = count_vect.transform(X_test['summary'])
20 #df4 = pd.DataFrame(final_test_summary.toarray(), columns=tf_idf_vect.get_fea
21
22 #f1_test_sparse = scipy.sparse.coo_matrix(df4['abandon'])
23 #f1_test_sparse = f1_test_sparse.transpose()
24 final_test_count = scipy.sparse.hstack([final_test_count, final_test_summary[
25
26
27 #f1_sparse = scipy.sparse.coo_matrix(df3['zucchini'])
28 #f1_sparse = f1_sparse.transpose()
29 final_tf_idf = scipy.sparse.hstack([final_tf_idf, final_summary[:,112]])
30 print("final_counts_tr_cv.shape = ",final_tf_idf.shape)
31
32 #final_test_summary = count_vect.transform(X_test['summary'])
33 #df4 = pd.DataFrame(final_test_summary.toarray(), columns=tf_idf_vect.get_fea
34
35 #f1_test_sparse = scipy.sparse.coo_matrix(df4['zucchini'])
36 #f1_test_sparse = f1_test_sparse.transpose()
37 final_test_count = scipy.sparse.hstack([final_test_count, final_test_summary[
38 print("test.shape = ",final_test_count.shape)
39
40 con_mat=nb_results(10,final_tf_idf,final_test_count,y_1,y_test)

```

```
final_counts_tr_cv.shape = (61441, 46448)
```

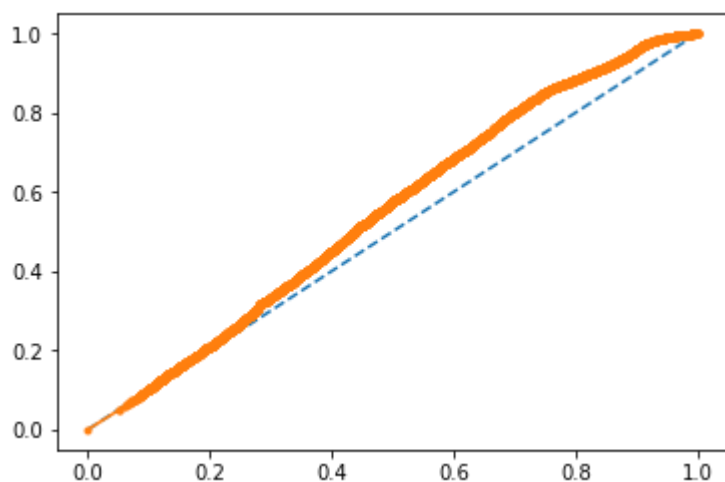
```
final_counts_tr_cv.shape = (61441, 46449)
```

```
test.shape = (26332, 46449)
```

The accuracy of the NB classifier for alpha = 10.000000 is 83.947288%

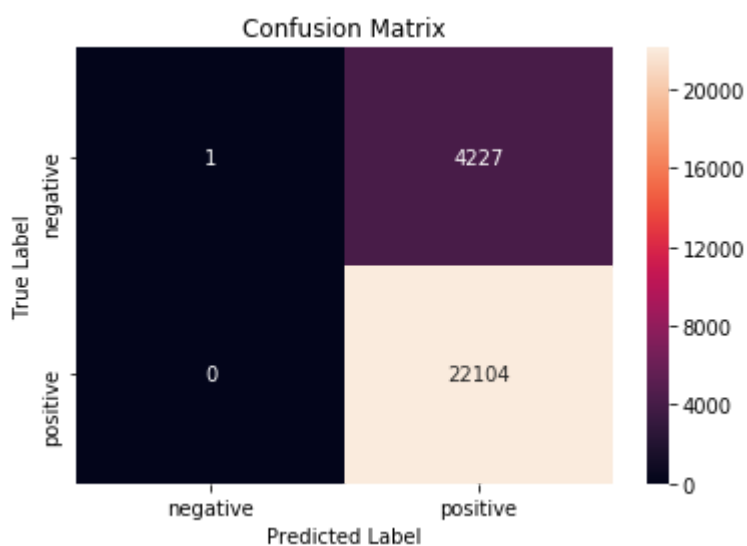
AUC: 0.550





Observation : With alpha = 10 the accuracy is 84 % and AUC is 0.550. This model is not good.

In [201]: 1 showHeatMap(con\_mat)



Observation: My model predicted 4227 points wrongly

## [6] Conclusions

Method	No of samples	alpha	accuracy	AUC Score
BOW	100000	10	84	0.738
TFIDF	100000	10	83	0.748
BOW1	100000	10	84	0.745
TFIDF1	100000	10	84	0.550

In [ ]: 1

