# Amazon Fine Food Reviews Analysis

Data Source: https://www.kaggle.com/snap/amazon-fine-food-reviews

EDA: https://nycdatascience.com/blog/student-works/amazon-fine-foods-visualization/

The Amazon Fine Food Reviews dataset consists of reviews of fine foods from Amazon.

Number of reviews: 568,454
Number of users: 256,059
Number of products: 74,258
Timespan: Oct 1999 - Oct 2012
Number of Attributes/Columns in data: 10

Attribute Information:

1. Id
2. ProductId - unique identifier for the product
3. UserId - unqiue identifier for the user
4. ProfileName
5. HelpfulnessNumerator - number of users who found the review helpful
6. HelpfulnessDenominator - number of users who indicated whether they found the review helpful or not
7. Score - rating between 1 and 5
8. Time - timestamp for the review
9. Summary - brief summary of the review
10. Text - text of the review

**Objective:**

Given a review, determine whether the review is positive (rating of 4 or 5) or negative (rating of 1 or 2).

[Q] How to determine if a review is positive or negative?

[Ans] We could use Score/Rating. A rating of 4 or 5 can be cosnidered as a positive review. A rating of 1 or 2 can be considered as negative one. A review of rating 3 is considered nuetral and such reviews are ignored from our analysis. This is an approximate and proxy way of determining the polarity (positivity/negativity) of a review.

# [1]. Reading Data

## [1.1] Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

In order to load the data, We have used the SQLITE dataset as it is easier to query the data and visualise the data efficiently.

Here as we only want to get the global sentiment of the recommendations (positive or negative), we will purposefully ignore all Scores equal to 3. If the score is above 3, then the recommendation wil be set to "positive". Otherwise, it will be set to "negative".

In [256]:
```python
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")


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

```python
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer

import re
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
import string
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer

from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle

from tqdm import tqdm
import os

from sklearn.cross_validation import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.cross_validation import cross_val_score
from collections import Counter
from sklearn.metrics import accuracy_score
from sklearn import cross_validation
```

In [257]:
```python
# using SQLite Table to read data.
con = sqlite3.connect('/Users/puravshah/Downloads/amazon-fine-food-reviews/database.sqlite')

# filtering only positive and negative reviews i.e.
# not taking into consideration those reviews with Score=3
```

```python
# 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 Sco
re != 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 sc
ore<3 a negative rating(0).
def partition(x):
    if x < 3:
        return 0
    return 1

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

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

Out[257]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfulnes |
|---|---|---|---|---|---|---|
| 0 | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1 | 1 |

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfules |
|---|---|---|---|---|---|---|
| 1 | 2 | B00813GRG4 | A1D87F6ZCVE5NK | dll pa | 0 | 0 |
| 2 | 3 | B000LQOCH0 | ABXLMWJIXXAIN | Natalia Corres "Natalia Corres" | 1 | 1 |

In [258]:
```python
display = pd.read_sql_query("""
SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
FROM Reviews
GROUP BY UserId
HAVING COUNT(*)>1
""", con)
```

In [259]:
```python
print(display.shape)
display.head()
```

(80668, 7)

Out[259]:

| | UserId | ProductId | ProfileName | Time | Score | Text | COUN |
|---|---|---|---|---|---|---|---|

| | UserId | ProductId | ProfileName | Time | Score | Text | COUI |
|---|---|---|---|---|---|---|---|
| 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 [260]: display[display['UserId']=='AZY10LLTJ71NX']
```

Out[260]:

| | UserId | ProductId | ProfileName | Time | Score | Text | C |
|---|---|---|---|---|---|---|---|

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

```
In [261]: display['COUNT(*)'].sum()
```

Out[261]: 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 [262]: display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID
""", con)
display.head()
```

Out[262]:

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

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfuln |
|---|---|---|---|---|---|---|
| 0 | 78445 | B000HDL1RQ | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |
| 1 | 138317 | B000HDOPYC | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |
| 2 | 138277 | B000HDOPYM | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |
| 3 | 73791 | B000HDOPZG | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |
| 4 | 155049 | B000PAQ75C | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |

As it can be seen above that same user has multiple reviews with same values for HelpfulnessNumerator, HelpfulnessDenominator, Score, Time, Summary and Text and on doing analysis it was found that

ProductId=B000HDOPZG was Loacker Quadratini Vanilla Wafer Cookies, 8.82-Ounce Packages (Pack of 8)

ProductId=B000HDL1RQ was Loacker Quadratini Lemon Wafer Cookies, 8.82-Ounce Packages (Pack of 8) and so on

It was inferred after analysis that reviews with same parameters other than ProductId belonged to the same product just having different flavour or quantity. Hence in order to reduce redundancy it was decided to eliminate the rows having same parameters.

The method used for the same was that we first sort the data according to ProductId and then just keep the first similar product review and delelte the others. for eg. in the above just the review for ProductId=B000HDL1RQ remains. This method ensures that there is only one representative for each product and deduplication without sorting would lead to possibility of different representatives still existing for the same product.

```python
In [263]: #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')
```

```python
In [264]: #Deduplication of entries
          final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"}, keep='first', inplace=False)
          final.shape
```

```
Out[264]: (87775, 10)
```

```python
In [265]: #Checking to see how much % of data still remains
          (final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

```
Out[265]: 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 calcualtions

```
In [266]: display= pd.read_sql_query("""
          SELECT *
          FROM Reviews
          WHERE Score != 3 AND Id=44737 OR Id=64422
          ORDER BY ProductID
          """, con)

          display.head()
```

Out[266]:

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

```
In [267]: final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [268]: #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[268]: 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 [269]:
```
# printing some random reviews
sent_0 = final['Text'].values[0]
print(sent_0)
print("="*50)
```

```python
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 [270]:
```python
# remove urls from text python: https://stackoverflow.com/a/40823105/40
84039
sent_0 = re.sub(r"http\S+", "", sent_0)
sent_1000 = re.sub(r"http\S+", "", sent_1000)
```

```
sent_150 = re.sub(r"http\S+", "", sent_1500)
sent_4900 = re.sub(r"http\S+", "", sent_4900)

print(sent_0)
```

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 [271]:
```
# 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 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 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 [272]:
```python
# https://stackoverflow.com/a/47091490/4084039
import re

def decontracted(phrase):
    # specific
    phrase = re.sub(r"won't", "will not", phrase)
    phrase = re.sub(r"can\'t", "can not", phrase)

    # general
    phrase = re.sub(r"n\'t", " not", phrase)
    phrase = re.sub(r"\'re", " are", phrase)
    phrase = re.sub(r"\'s", " is", phrase)
    phrase = re.sub(r"\'d", " would", phrase)
    phrase = re.sub(r"\'ll", " will", phrase)
    phrase = re.sub(r"\'t", " not", phrase)
    phrase = re.sub(r"\'ve", " have", phrase)
    phrase = re.sub(r"\'m", " am", phrase)
    return phrase
```

In [273]:
```python
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 [274]: 
```python
#remove words with numbers python: https://stackoverflow.com/a/1808237
0/4084039
sent_0 = re.sub("\S*\d\S*", "", sent_0).strip()
print(sent_0)
```

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 [275]: 
```python
#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 [276]: 
```python
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'no
t'
# <br /><br /> ==> after the above steps, we are getting "br br"
# we are including them into stop words list
# instead of <br /> if we have <br/> these tags would have revmoved in
 the 1st step

stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'o
urs', 'ourselves', 'you', "you're", "you've",\
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselve
s', 'he', 'him', 'his', 'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'it
s', 'itself', 'they', 'them', 'their',\
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'th
is', 'that', "that'll", 'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'h
ave', 'has', 'had', 'having', 'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or',
 'because', 'as', 'until', 'while', 'of', \
```

```
                    'at', 'by', 'for', 'with', 'about', 'against', 'between',
'into', 'through', 'during', 'before', 'after',\
                    'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out',
'on', 'off', 'over', 'under', 'again', 'further',\
                    'then', 'once', 'here', 'there', 'when', 'where', 'why', 'h
ow', 'all', 'any', 'both', 'each', 'few', 'more',\
                    'most', 'other', 'some', 'such', 'only', 'own', 'same', 's
o', 'than', 'too', 'very', \
                    's', 't', 'can', 'will', 'just', 'don', "don't", 'should',
"should've", 'now', 'd', 'll', 'm', 'o', 're', \
                    've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't",
'didn', "didn't", 'doesn', "doesn't", 'hadn',\
                    "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "is
n't", 'ma', 'mightn', "mightn't", 'mustn',\
                    "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn',
 "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
                    'won', "won't", 'wouldn', "wouldn't"])
```

In [277]:
```python
# Combining all the above stundents
from tqdm import tqdm
preprocessed_reviews = []
# tqdm is for printing the status bar
for sentance in tqdm(final['Text'].values):
    sentance = re.sub(r"http\S+", "", sentance)
    sentance = BeautifulSoup(sentance, 'lxml').get_text()
    sentance = decontracted(sentance)
    sentance = re.sub("\S*\d\S*", "", sentance).strip()
    sentance = re.sub('[^A-Za-z]+', ' ', sentance)
    # https://gist.github.com/sebleier/554280
    sentance = ' '.join(e.lower() for e in sentance.split() if e.lower
() not in stopwords)
    preprocessed_reviews.append(sentance.strip())
```

```
100%|████████████| 87773/87773 [00:25<00:00, 3422.82it/s]
```

In [278]:
```python
preprocessed_reviews[1500]
```

Out[278]: `'way hot blood took bite jig lol'`

## <span style="color:red">[3.2] Preprocessing Review Summary</span>

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

# [4] Featurization

## [4.1] BAG OF WORDS

In [311]:
```python
X_1, X_test, y_1, y_test = cross_validation.train_test_split(preprocess
ed_reviews, final['Score'], test_size=0.3, random_state=0)
X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_1, y_1, te
st_size=0.3)
count_vect = CountVectorizer() #in scikit-learn
count_vect.fit(X_tr)
count_vect.fit(X_cv)
count_vect.fit(X_test)
print("some feature names ", count_vect.get_feature_names()[:10])
print('='*50)

final_counts = count_vect.transform(X_tr)
X_cv_bow=count_vect.transform(X_cv)
X_test_bow=count_vect.transform(X_test)
print("the type of count vectorizer ",type(final_counts))
print("the shape of out text BOW vectorizer ",final_counts.get_shape())
print("the number of unique words ", final_counts.get_shape()[1])
```
```
some feature names  ['aa', 'aaa', 'aaaa', 'aaaand', 'aaah', 'aaahs', 'a
afco', 'aahhhs', 'aahing', 'aamazon']
==================================================
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (43008, 30512)
the number of unique words  30512
```

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

In [280]:
```python
#bi-gram, tri-gram and n-gram

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

# you can choose these numebrs min_df=10, max_features=5000, of your ch
oice
count_vect = CountVectorizer(ngram_range=(1,2), min_df=10, max_features
=5000)
final_bigram_counts = count_vect.fit_transform(X_tr)
X_cv_ngram = count_vect.fit_transform(X_cv)
X_test_ngram=count_vect.fit_transform(X_test)
print("the type of count vectorizer ",type(final_bigram_counts))
print("the shape of out text BOW vectorizer ",final_bigram_counts.get_s
hape())
print("the number of unique words including both unigrams and bigrams "
, final_bigram_counts.get_shape()[1])
```

```
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (43008, 5000)
the number of unique words including both unigrams and bigrams  5000
```

## [4.3] TF-IDF

In [281]:
```python
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
tf_idf_vect.fit(X_tr)
tf_idf_vect.fit(X_cv)
tf_idf_vect.fit(X_test)
print("some sample features(unique words in the corpus)",tf_idf_vect.ge
t_feature_names()[0:10])
print('='*50)
```

```
final_tf_idf = tf_idf_vect.transform(X_tr)
X_cv_tfidf=tf_idf_vect.transform(X_cv)
X_test_tfidf=tf_idf_vect.transform(X_test)
print("the type of count vectorizer ",type(final_tf_idf))
print("the shape of out text TFIDF vectorizer ",final_tf_idf.get_shape
())
print("the number of unique words including both unigrams and bigrams "
, final_tf_idf.get_shape()[1])
```

```
some sample features(unique words in the corpus) ['ability', 'able', 'a
ble buy', 'able drink', 'able eat', 'able enjoy', 'able find', 'able ge
t', 'able give', 'able make']
==================================================
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer  (43008, 15763)
the number of unique words including both unigrams and bigrams  15763
```

## [4.4] Word2Vec

In [282]:
```python
# Train your own Word2Vec model using your own text corpus
i=0
list_of_sentance=[]
for sentance in preprocessed_reviews:
    list_of_sentance.append(sentance.split())
```

In [283]:
```python
# Using Google News Word2Vectors

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

```python
# it's 1.9GB in size.


# http://kavita-ganesan.com/gensim-word2vec-tutorial-starter-code/#.W17
SRFAzZPY
# you can comment this whole cell
# or change these varible according to your need

is_your_ram_gt_16g=False
want_to_use_google_w2v = False
want_to_train_w2v = True

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

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

```
[('fantastic', 0.8536819815635681), ('awesome', 0.8350322842597961),
('excellent', 0.8242692351341248), ('good', 0.8191254734992981), ('terr
ific', 0.7985398173332214), ('wonderful', 0.7982454299926758), ('perfec
t', 0.770243763923645), ('amazing', 0.7415543794631958), ('fabulous',
0.7378154993057251), ('nice', 0.7201124429702759)]
==================================================
[('greatest', 0.8203911781311035), ('tastiest', 0.7233695983886719),
('best', 0.7149415016174316), ('nastiest', 0.6551864147186279), ('ive',
0.6354540586471558), ('closest', 0.6261584758758545), ('surpass', 0.602
5659441947937), ('weakest', 0.6011222004890442), ('horrible', 0.5994516
611099243), ('awful', 0.5992503762245178)]
```

```
In [284]: w2v_words = list(w2v_model.wv.vocab)
          print("number of words that occured minimum 5 times ",len(w2v_words))
          print("sample words ", w2v_words[0:50])
```

```
number of words that occured minimum 5 times  17386
sample words  ['dogs', 'loves', 'chicken', 'product', 'china', 'wont',
'buying', 'anymore', 'hard', 'find', 'products', 'made', 'usa', 'one',
'isnt', 'bad', 'good', 'take', 'chances', 'till', 'know', 'going', 'imp
orts', 'love', 'saw', 'pet', 'store', 'tag', 'attached', 'regarding',
'satisfied', 'safe', 'infestation', 'literally', 'everywhere', 'flyin
g', 'around', 'kitchen', 'bought', 'hoping', 'least', 'get', 'rid', 'we
eks', 'fly', 'stuck', 'squishing', 'buggers', 'success', 'rate']
```

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

### [4.4.1.1] Avg W2v

```python
In [285]: # average Word2Vec
          # compute average word2vec for each review.
          sent_vectors = []; # the avg-w2v for each sentence/review is stored in
           this list
          for sent in tqdm(list_of_sentance): # for each review/sentence
              sent_vec = np.zeros(50) # as word vectors are of zero length 50, yo
          u might need to change this to 300 if you use google's w2v
              cnt_words =0; # num of words with a valid vector in the sentence/re
          view
              for word in sent: # for each word in a review/sentence
                  if word in w2v_words:
                      vec = w2v_model.wv[word]
                      sent_vec += vec
                      cnt_words += 1
              if cnt_words != 0:
                  sent_vec /= cnt_words
              sent_vectors.append(sent_vec)
```

```
print(len(sent_vectors))
print(len(sent_vectors[0]))
```

```
100%|████████| 87773/87773 [02:15<00:00, 648.11it/s]
```

```
87773
50
```

**[4.4.1.2] TFIDF weighted W2v**

In [286]:
```python
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer()
tf_idf_matrix = model.fit_transform(preprocessed_reviews)
# we are converting a dictionary with word as a key, and the idf as a v
alue
dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
```

In [287]:
```python
# TF-IDF weighted Word2Vec
tfidf_feat = model.get_feature_names() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and ce
ll_val = tfidf

tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is st
ored in this list
row=0;
for sent in tqdm(list_of_sentance): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum =0; # num of words with a valid vector in the sentence/r
eview
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
#            tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
```

```
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    tfidf_sent_vectors.append(sent_vec)
    row += 1
```
100%|████████████| 87773/87773 [24:15<00:00, 56.02it/s]

# [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 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 values of `feature_log_prob_` parameter of MultinomialNB 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 with predicted and original labels of test data points. Please visualize your confusion matrices using seaborn heatmaps.


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



**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 Multinomial Naive Bayes

## [5.1] Applying Naive Bayes on BOW, <span style="color:red">SET 1</span>

In [323]:
```python
#Some parts of code taken from the reference ipynb provided

from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import roc_auc_score
auc=[]
auc2=[]
#defining alpha values
alpha=[0.00001,0.0001,0.001,0.01,0.1,1,10,100,1000,10000]
print(np.log10(alpha))
#for different alpha values finding the auc values to be plotted
for i in alpha:
    model=MultinomialNB(alpha=i)
    model.fit(final_counts,y_tr)
    pred=model.predict(X_cv_bow)
    probab=model.predict_proba(X_cv_bow)
    fpr,tpr,thresholds=metrics.roc_curve(y_cv,pred)
    auc.append(roc_auc_score(y_cv,pred))
    acc=accuracy_score(y_cv,pred,normalize=True)*float(100)
    pred2=model.predict(final_counts)
    probab2=model.predict_proba(final_counts)
    fpr_train,tpr_train,thresholds=roc_curve(y_tr,pred2)
    auc2.append(roc_auc_score(y_tr,pred2))
    print('Accuracy score at alpha value %d is %d'%(i,acc))
plt.plot(np.log10(alpha),auc,'b')
plt.plot(np.log10(alpha),auc2,'y')
plt.title('Auc curve for different values of alpha')
plt.xlabel('Alpha values')
plt.ylabel('AUC')

#it can be seen from the plot that the auc values for the train and cv
 dataset overlap and there is no clear distinction between the two
```
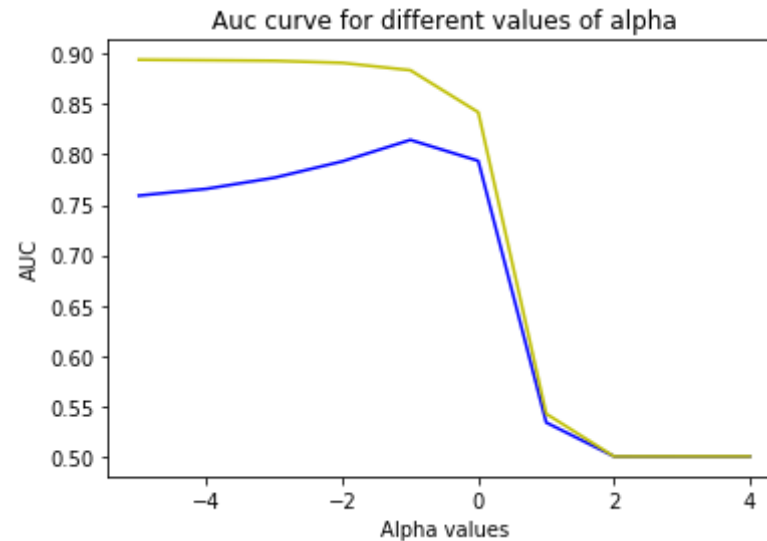
```
[-5. -4. -3. -2. -1.  0.  1.  2.  3.  4.]
Accuracy score at alpha value 0 is 88
```

```
Accuracy score at alpha value 0 is 88
Accuracy score at alpha value 0 is 88
Accuracy score at alpha value 0 is 89
Accuracy score at alpha value 0 is 89
Accuracy score at alpha value 0 is 90
Accuracy score at alpha value 1 is 90
Accuracy score at alpha value 10 is 84
Accuracy score at alpha value 100 is 83
Accuracy score at alpha value 1000 is 83
Accuracy score at alpha value 10000 is 83
```
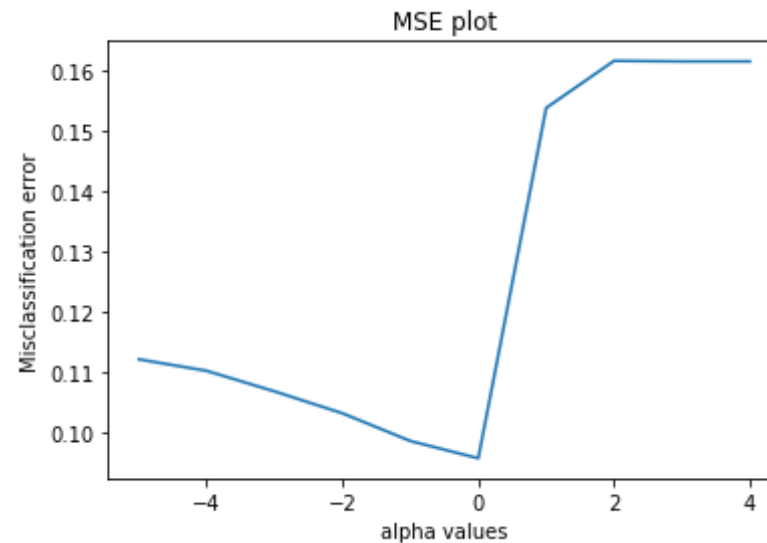
Out[323]: `Text(0,0.5,'AUC')`



In [324]:
```python
#This part of the code is to caluclate the optimal alpha by calculating
 the mis-classification error
cv_scores=[]
alpha=[0.00001,0.0001,0.001,0.01,0.1,1,10,100,1000,10000]
for i in alpha:
    model=MultinomialNB(alpha=i)
    scores=cross_val_score(model,final_counts,y_tr,cv=10,scoring='accur
acy')
    cv_scores.append(scores.mean())
MSE=[1-x for x in cv_scores]
```

```python
plt.plot(np.log10(alpha),MSE)
plt.title('MSE plot')
plt.xlabel('alpha values')
plt.ylabel('Misclassification error')
optimal_alpha=alpha[MSE.index(min(MSE))]
print('The optimal alpha vale for low error is=%d'%optimal_alpha)
```
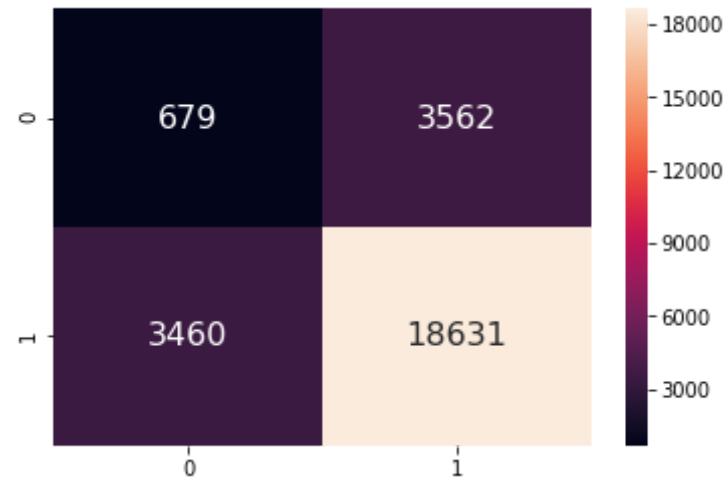
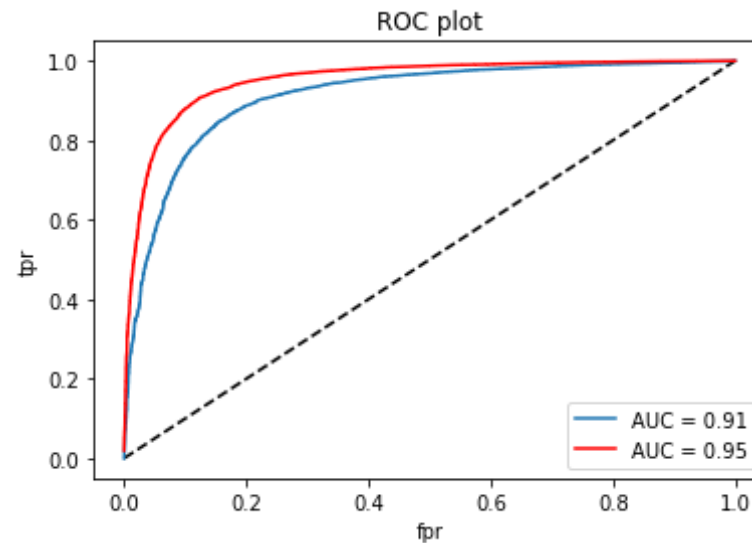The optimal alpha vale for low error is=1



```python
In [325]: #Using the optimal value of alpha that was calculated we run our model
           with that alpha value and check the accuracy
          model=MultinomialNB(alpha=optimal_alpha)
          model.fit(final_counts,y_tr)
          predict=model.predict(X_test_bow)
          probab=model.predict_proba(X_test_bow)
          probab=probab[:,1]
          predict2=model.predict(final_counts)
          probab2=model.predict_proba(final_counts)
          probab2=probab2[:,1]
          acc=accuracy_score(y_test,predict,normalize=True)*float(100)
          print('The accuracy of the model with optimal value of alpha is=%d'%opt
          imal_alpha)
```

The accuracy of the model with optimal value of alpha is=1

In [326]:
```python
from sklearn.metrics import confusion_matrix as cm
from sklearn.metrics import roc_auc_score
#plotting the confusion matrix and the ROC curve
df_cm = pd.DataFrame(confusion_matrix(Y_test, predict), range(2),range(
2))
    #heatman for visualization of matrix
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
plt.show()
fpr,tpr,thresholds=roc_curve(y_test,probab)
auc=roc_auc_score(y_test,probab)
fpr_train,tpr_train,thresholds=roc_curve(y_tr,probab2)
auc2=roc_auc_score(y_tr,probab2)
plt.plot([0,1],[0,1],'k--')
plt.plot(fpr,tpr,label = 'AUC = %0.2f' %auc)
plt.plot(fpr_train,tpr_train,'r',label = 'AUC = %0.2f' %auc2)
plt.title('ROC plot')
plt.xlabel('fpr')
plt.ylabel('tpr')
plt.legend(loc='lower right')
plt.show()
```

ROC plot

**[5.1.1] Top 10 important features of positive class from <span style="color:red">SET 1</span>**

In [327]:
```python
#calculating the feature log probabilities
feature_log_prob=model.feature_log_prob_
#print(feature_log_prob.shape)
#concerting it into a data frame
df=pd.DataFrame(feature_log_prob)
#print(df.head())
#Getting the feature names from bow vectorizer
bow=count_vect.get_feature_names()
#Making the feature names as the index(matching the probabilities to th
e respective values)
df=pd.DataFrame(feature_log_prob,columns=bow)
#print(df.head())
df=df.T
#print(df.head())
#printing the top 10 positive features
print("Top 10 Positive Features:-\n",df[1].sort_values(ascending = Fals
e)[0:10])
```

```
Top 10 Positive Features:-
 not       -3.720529
like      -4.530450
good      -4.668716
great     -4.755367
one       -4.895308
taste     -4.959547
coffee    -5.006954
flavor    -5.063586
love      -5.072091
would     -5.076085
Name: 1, dtype: float64
```

**[5.1.2] Top 10 important features of negative class from <span style="color:red">SET 1</span>**

In [328]:
```python
# printing the top 10 negative features
print("Top 10 negative Features:-\n",df[0].sort_values(ascending = False)[0:10])
```

```
Top 10 negative Features:-
 not        -3.343000
like       -4.469759
would      -4.732575
taste      -4.759215
product    -4.781486
one        -4.960399
coffee     -5.164503
good       -5.214759
flavor     -5.232849
no         -5.259440
Name: 0, dtype: float64
```

## [5.2] Applying Naive Bayes on TFIDF, <span style="color:red">SET 2</span>

In [341]:
```python
#The below code is exactly same as the previous one with the only difference being that it is applied to TF-IDF vectorizer
```

```python
auc=[]
auc2=[]
alpha=[0.00001,0.0001,0.001,0.01,0.1,1,10,100,1000,10000]
for i in alpha:
    model=MultinomialNB(alpha=i)
    model.fit(final_tf_idf,y_tr)
    pred1=model.predict(X_cv_tfidf)
    probab1=model.predict_proba(X_cv_tfidf)
    probab1=probab1[:,1]
    fpr,tpr,thresholds=metrics.roc_curve(y_cv,probab1)
    auc.append(roc_auc_score(y_cv,probab1))
    acc=accuracy_score(y_cv,pred1,normalize=True)*float(100)
    pred4=model.predict(final_tf_idf)
    probab4=model.predict_proba(final_tf_idf)
    probab4=probab4[:,1]
    fpr_train,tpr_train,thresholds=metrics.roc_curve(y_tr,probab4)
    auc2.append(roc_auc_score(y_tr,probab4))
    print('Accuracy score at alpha value %d is %d'%(i,acc))
plt.plot(np.log10(alpha),auc)
plt.plot(np.log10(alpha),auc2,'r')
plt.title('Auc curve for different values of alpha')
plt.xlabel('alpha')
plt.ylabel('auc')
```
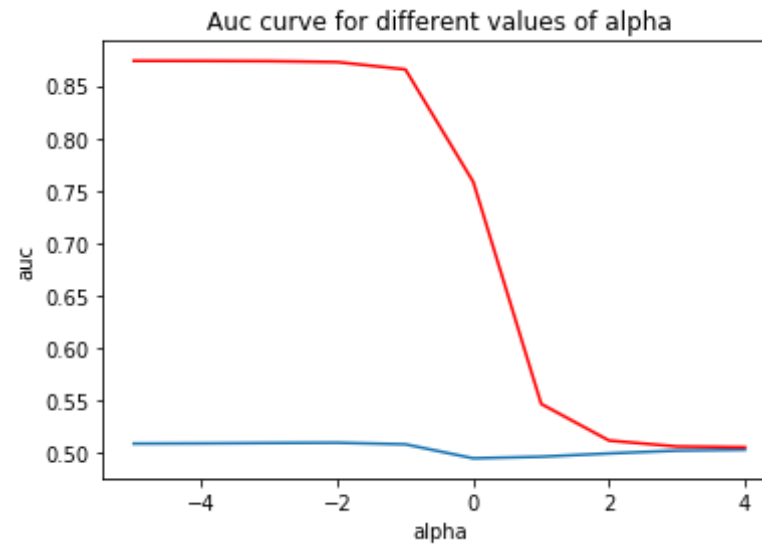
```
Accuracy score at alpha value 0 is 83
Accuracy score at alpha value 0 is 83
Accuracy score at alpha value 0 is 83
Accuracy score at alpha value 0 is 83
Accuracy score at alpha value 0 is 83
Accuracy score at alpha value 1 is 83
Accuracy score at alpha value 10 is 83
Accuracy score at alpha value 100 is 83
Accuracy score at alpha value 1000 is 83
Accuracy score at alpha value 10000 is 83
```
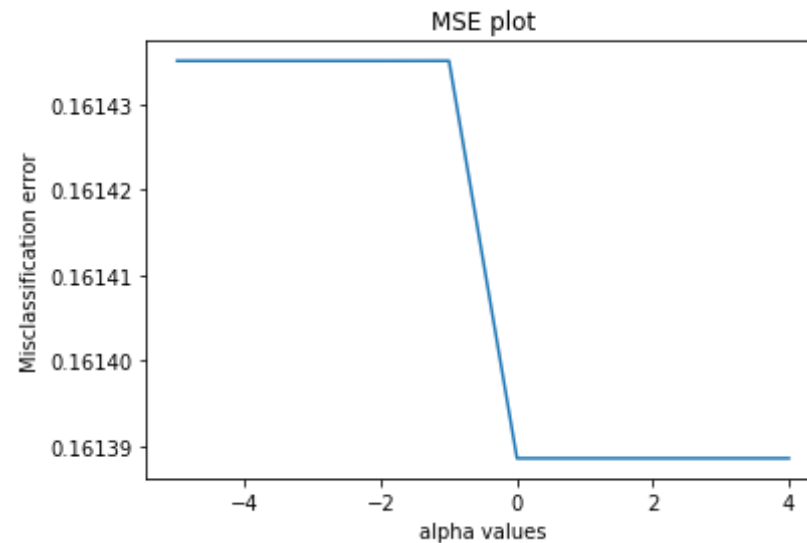
Out[341]: Text(0,0.5,'auc')

Auc curve for different values of alpha



In [342]:
```python
# Please write all the code with proper documentation
cv_scores=[]
alpha=[0.00001,0.0001,0.001,0.01,0.1,1,10,100,1000,10000]
for i in alpha:
    model=MultinomialNB(alpha=i)
    scores=cross_val_score(model,final_tf_idf,y_tr,cv=10,scoring='accuracy')
    cv_scores.append(scores.mean())
MSE=[1-x for x in cv_scores]
plt.plot(np.log10(alpha),MSE)
plt.title('MSE plot')
plt.xlabel('alpha values')
plt.ylabel('Misclassification error')
optimal_alpha=alpha[MSE.index(min(MSE))]
print('The optimal alpha vale for low error is=%d'%optimal_alpha)
```
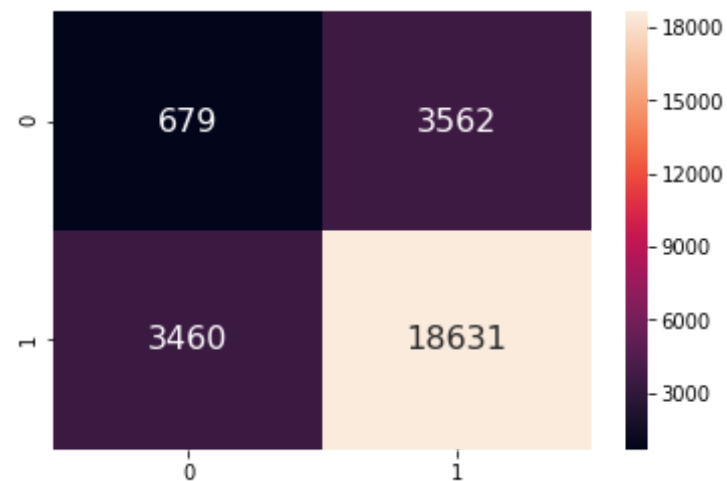
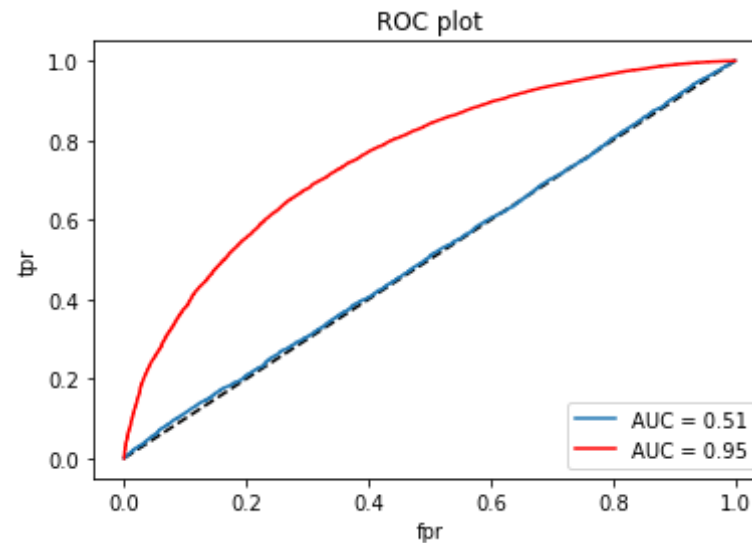The optimal alpha vale for low error is=1

MSE plot

In [343]:
```python
model=MultinomialNB(alpha=optimal_alpha)
model.fit(final_tf_idf,y_tr)
predict1=model.predict(X_test_tfidf)
probab1=model.predict_proba(X_test_tfidf)
probab1=probab1[:,1]
predict3=model.predict(final_tf_idf)
probab3=model.predict_proba(final_tf_idf)
probab3=probab3[:,1]
acc=accuracy_score(y_test,predict1,normalize=True)*float(100)
print('The accuracy of the model with optimal value of alpha is=%d'%opt
imal_alpha)
```

The accuracy of the model with optimal value of alpha is=1

In [344]:
```python
df_cm = pd.DataFrame(confusion_matrix(Y_test, predict), range(2),range(
2))
    #heatman for visualization of matrix
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
plt.show()
fpr,tpr,thresholds=roc_curve(y_test,probab1)
auc1=roc_auc_score(y_test,probab1)
```

```
fpr_train,tpr_train,thresholds=roc_curve(y_tr,probab3)
auc3=roc_auc_score(y_tr,probab2)
plt.plot([0,1],[0,1],'k--')
plt.plot(fpr,tpr,label = 'AUC = %0.2f' %auc1)
plt.plot(fpr_train,tpr_train,'r',label = 'AUC = %0.2f' %auc3)
plt.title('ROC plot')
plt.xlabel('fpr')
plt.ylabel('tpr')
plt.legend(loc='lower right')
plt.show()
```

ROC plot

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

```python
feature_log_prob_tfidf=model.feature_log_prob_
#print(feature_log_prob.shape)
df=pd.DataFrame(feature_log_prob_tfidf)
#print(df.shape)
tfidf_features=tf_idf_vect.get_feature_names()
df1=pd.DataFrame(feature_log_prob_tfidf,columns=tfidf_features)
#print(df1.shape)
#print(df.head())
df1=df1.T
#print(df1.head())
print("Top 10 Positive Features:-\n",df1[1].sort_values(ascending = False)[0:10])
```

```
Top 10 Positive Features:-
 not      -4.999648
like      -5.551137
good      -5.583325
great     -5.597252
```

```
coffee    -5.624187
taste     -5.753221
product   -5.753313
tea       -5.771656
one       -5.809017
love      -5.812956
Name: 1, dtype: float64
```

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

In [339]:
```python
# Please write all the code with proper documentation
print("Top 10 negative Features:-\n",df1[0].sort_values(ascending = False)[0:10])
```

```
Top 10 negative Features:-
 not        -5.016674
like       -5.555650
great      -5.564681
good       -5.579038
coffee     -5.596584
product    -5.727963
taste      -5.751141
tea        -5.760627
love       -5.788780
one        -5.820209
Name: 0, dtype: float64
```

# [6] Conclusions

In [340]:
```python
# Please compare all your models using Prettytable library
from prettytable import PrettyTable
x=PrettyTable()
x.field_names=['Vectorizer','Model','Hyperparameter','AUC']
x.add_row(["Bag Of Words","Multinomial Naive bayes",1,0.91])
x.add_row(["TF-IDF","Multinomial Naive Bayes",0.01,0.48])
print(x)
```

```
+--------------+--------------------------+----------------+------+
|  Vectorizer  |          Model           | Hyperparameter | AUC  |
+--------------+--------------------------+----------------+------+
| Bag Of Words | Multinomial Naive bayes  |       1        | 0.91 |
|    TF-IDF    | Multinomial Naive Bayes  |      0.01      | 0.95 |
+--------------+--------------------------+----------------+------+
```