

Amazon Fine Food Reviews Analysis

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

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

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

Number of reviews: 568,454

Number of users: 256,059

Number of products: 74,258

Timespan: Oct 1999 - Oct 2012

Number of Attributes/Columns in data: 10

Attribute Information:

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

Objective:

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

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

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

[1]. Reading Data

[1.1] Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

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

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

In [1]:

```
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
```

```

import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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

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

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

from tqdm import tqdm
import os

```

```

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

```

In [3]:

```

# using SQLite Table to read data.
con = sqlite3.connect('database.sqlite')

# filtering only positive and negative reviews i.e.
# not taking into consideration those reviews with Score=3
# SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000, will give top 500000 data points
# you can change the number to any other number based on your computing power

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

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

# Give reviews with Score>3 a positive rating, and reviews with a score<3 a negative rating.
def partition(x):
    if x < 3:
        return 'negative'
    return 'positive'

#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)
print("Data Points in Each class :")
print(filtered_data['Score'].value_counts())
filtered_data.head(3)

```

```

Number of data points in our data (525814, 10)
Data Points in Each class :
positive      443777
negative      82037
Name: Score, dtype: int64

```

Out[3]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | Ti |
|---|----|------------|----------------|-------------|----------------------|------------------------|----------|----------|
| 0 | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1 | 1 | positive | 13038624 |

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | Ti |
|---|----|------------|----------------|--|----------------------|------------------------|----------|----------|
| 1 | 2 | B00813GRG4 | A1D87F6ZCVE5NK | dll pa | 0 | 0 | negative | 13469760 |
| 2 | 3 | B000LQOCH0 | ABXLMWJIXXAIN | Natalia Corres "Natalia Corres" | 1 | 1 | positive | 12190176 |

[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 [4]:

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

Out[4]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | Ti |
|---|--------|------------|---------------|--------------------|----------------------|------------------------|-------|----------|
| 0 | 78445 | B000HDL1RQ | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 11995776 |
| 1 | 138317 | B000HDOPYC | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 11995776 |
| 2 | 138277 | B000HDOPYM | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 11995776 |
| 3 | 73791 | B000HDOPZG | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 11995776 |
| 4 | 155049 | B000PAQ75C | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 11995776 |

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 [5]:

```
#Sorting the data taking productid as the parameter
sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True, inplace=False, kind='quicksort', na_position='last')
sorted_data.shape
```

Out[5]:

(525814, 10)

In [6]:

```
#Deleting the duplicates reviews which is created when user writed a review for the product, it automatically generates for the same product of different color etc
final = sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"}, keep='first', inplace=False)
final.shape
```

Out[6]:

(364173, 10)

In [7]:

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

Out[7]:

69.25890143662969

In [8]:

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

display.head()
```

Out[8]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | Time |
|---|-------|------------|----------------|----------------------------|----------------------|------------------------|-------|----------|
| 0 | 64422 | B000MIDROQ | A161DK06JJMCYF | J. E. Stephens "Jeanne" | 3 | 1 | 5 | 12248926 |

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | Time |
|---|-------|------------|----------------|-------------|----------------------|------------------------|-------|----------|
| 1 | 44737 | B001EQ55RW | A2V0I904FH7ABY | Ram | 3 | 2 | 4 | 12128832 |

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 [9]:

```
#Dropping the data which has HelpfulnessNumerator<HelpfulnessDenominator which is impossible
final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
#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()
```

(364171, 10)

Out[9]:

```
positive    307061
negative     57110
Name: Score, dtype: int64
```

In [10]:

```
#Checking to see how much % of data still remains
print("Percentage of data still remains", (final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100)

print("Final Data", final.shape)
```

Percentage of data still remains 69.25852107399194
Final Data (364171, 10)

[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 [2]:

```
stop = set(stopwords.words('english')) #set of stopwords
sno = nltk.stem.SnowballStemmer('english') #initialising the snowball stemmer

def cleanhtml(sentence): #function to clean the word of any html-tags
```

```

cleanr = re.compile('<.*?>')
cleantext = re.sub(cleanr, ' ', sentence)
return cleantext
def cleanpunc(sentence): #function to clean the word of any punctuation or special characters
cleaned = re.sub(r'[?!|\\'|"|#]',r'',sentence)
cleaned = re.sub(r'[,|,|)(|\\|/]',r'',cleaned)
return cleaned

```

In [12]:

```

#Code for implementing step-by-step the checks mentioned in the pre-processing phase
# this code takes a while to run as it needs to run on 500k sentences.
from tqdm import tqdm
i=0
str1=' '
final_string=[]
all_positive_words=[] # store words from +ve reviews here
all_negative_words=[] # store words from -ve reviews here.
s=''
for sent in tqdm(final['Text'].values):
    filtered_sentence=[]
    #print(sent);
    sent=cleanhtml(sent) # remove HTML tags
    for w in sent.split():
        for cleaned_words in cleanpunc(w).split():
            if((cleaned_words.isalpha()) & (len(cleaned_words)>2)):
                if(cleaned_words.lower() not in stop):
                    s=(sno.stem(cleaned_words.lower()).encode('utf8'))
                    filtered_sentence.append(s)
                    if (final['Score'].values)[i] == 'positive':
                        all_positive_words.append(s) #list of all words used to describe positive r
reviews
                    if(final['Score'].values)[i] == 'negative':
                        all_negative_words.append(s) #list of all words used to describe negative r
reviews reviews
                else:
                    continue
            else:
                continue
    #print(filtered_sentence)
    str1 = b" ".join(filtered_sentence) #final string of cleaned words
    #print("*****")

    final_string.append(str1)
    i+=1

```

100% | 364171/364171
[09:42<00:00, 625.43it/s]

In [13]:

```

final['CleanedText']=final_string #adding a column of CleanedText which displays the data after pr
e-processing of the review
final['CleanedText']=final['CleanedText'].str.decode("utf-8")
final.head(3)

```

Out[13]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | |
|--------|--------|------------|---------------|--------------------|----------------------|------------------------|----------|----|
| 138706 | 150524 | 0006641040 | ACITT7DI6IDDL | shari zychinski | 0 | 0 | positive | 93 |
| | | | | | | | | |

| 138688 | 150506 | 0006641040 | A2IW4PEEK03800 | Used | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | 11 |
|--------|--------|------------|----------------|------|--------------------------|----------------------|------------------------|----------|----|
| | | | | | | | | | |
| 138689 | 150507 | 0006641040 | A1S4A3IQ2MU7V4 | | sally sue "sally sue" | 1 | 1 | positive | 11 |

In [4]:

```
final_data=final.sort_values('Time', axis=0, ascending=True, inplace=False, kind='quicksort', na_position='last')
```

In [6]:

```
final = final_data.head(100000)
```

In [7]:

```
X_train_data = final[:60000]
X_test_data = final[60000:100000]
y_train = X_train_data['Score']
y_test = X_test_data['Score']
print("Data")
print(X_train_data.shape)
print(X_test_data.shape)
print("Label")
print(y_train.shape)
print(y_test.shape)
```

```
Data
(60000, 11)
(40000, 11)
Label
(60000,)
(40000,)
```

[3.2] Preprocessing Review Summary

In [17]:

```
## Similarly you can do preprocessing for review summary also.
```

[4] Featurization

[4.1] BAG OF WORDS

In [9]:

```
#BoW on Text
print("**Bow Vectorizer**")
print("="*50)
count_vect = CountVectorizer(min_df = 50)
X_train_BOW = count_vect.fit_transform(X_train_data['CleanedText'])
X_test_BOW = count_vect.transform(X_test_data['CleanedText'])
print(X_train_BOW.shape)
print(X_test_BOW.shape)
```

```
**Bow Vectorizer**
```

(60000, 2951)
(40000, 2951)

[4.2] Bi-Grams and n-Grams.

In [19]:

```
# #bi-gram, tri-gram and n-gram

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

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

[4.3] TF-IDF

In [10]:

```
#TFIDF on Text
print("***TFIDF Vectorizer**")
print("=="*50)
tf_idf_vect = TfidfVectorizer(min_df = 50)
X_train_tfidf = tf_idf_vect.fit_transform(X_train_data['CleanedText'])
X_test_tfidf = tf_idf_vect.transform(X_test_data['CleanedText'])
print(X_train_tfidf.shape)
print(X_test_tfidf.shape)
```

```
**TFIDF Vectorizer**
```

```
(60000, 2951)
(40000, 2951)
```

[4.4] Word2Vec

In [51]:

```
import gensim
i=0
list_of_sent_train=[]
for sent in tqdm(X_train_data['Text'].values):
    filtered_sentence=[]
    sent=cleanhtml(sent)
    for w in sent.split():
        for cleaned_words in cleanpunc(w).split():
            if(cleaned_words.isalpha()): # checking is the word is alphabet
                filtered_sentence.append(cleaned_words.lower()) # appending to the list
            else:
                continue
    list_of_sent_train.append(filtered_sentence)
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 60000/60000  
[00:16<00:00, 3655.77it/s]
```

In [52]:

```
import gensim
i=0
list_of_sent_test=[]
for sent in tqdm(X_test_data['Text+'].values):
```



```
(('delicious', 0.712010070000217),
 ('tasty', 0.7192407846450806),
 ('moist', 0.7166850566864014),
 ('dense', 0.7083577513694763),
 ('addictive', 0.6990361213684082),
 ('nutritious', 0.6965190172195435)])
```

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

[4.4.1.1] Avg W2v

In [58]:

```
#TRAIN
# average Word2Vec
# compute average word2vec for each review.
sent_vectors_train = [];
for sent in tqdm(list_of_sent_train):
    sent_vec = np.zeros(50)
    cnt_words = 0;
    for word in sent: #
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    sent_vectors_train.append(sent_vec)
print(len(sent_vectors_train))
print(len(sent_vectors_train[0]))
```

[illegible]

60000
50

In [59]:

```
#TEST
# average Word2Vec
# compute average word2vec for each review
sent_vectors_test = [];
for sent in tqdm(list_of_sent_test):
    sent_vec = np.zeros(50)
    cnt_words = 0;
    for word in sent: #
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    sent_vectors_test.append(sent_vec)
print(len(sent_vectors_test))
print(len(sent_vectors_test[0]))
```

[illegible]

40000
50

[4.4.1.2] TFIDF weighted W2v

In [60]:

```
tfidf_vect = TfidfVectorizer(min_df = 50)
train_tfidf_w2v = tfidf_vect.fit_transform(X_train_data["CleanedText"])
```

(60000, 2951)
(40000, 2951)

```
# TF-IDF weighted Word2Vec
tfidf_feat = tfidf_vect.get_feature_names() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val = tfidf

tfidf_sent_vectors_train = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(list_of_sent_train): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum =0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
            #
            tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole corpus
            # sent.count(word) = tf value of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    tfidf_sent_vectors_train.append(sent_vec)
    row += 1
```

```
# TF-IDF weighted Word2Vec
tfidf_feat = tfidf_vect.get_feature_names() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val = tfidf

tfidf_sent_vectors_test = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(list_of_sent_test): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum =0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
            #
            tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole corpus
            # sent.count(word) = tf value 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_test.append(sent_vec)
    row += 1
```

1. Apply Random Forests & GBDT 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)
- **SET 3:** Review text, preprocessed one converted into vectors using (AVG W2v)
- **SET 4:** Review text, preprocessed one converted into vectors using (TFIDF W2v)

2. The hyper parameter tuning (Consider two hyperparameters: `n_estimators` & `max_depth`)

- Find the best hyper parameter which will give the maximum [AUC](#) value
- Find the best hyper parameter 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

- Get top 20 important features and represent them in a word cloud. Do this for BOW & TFIDF.

4. Feature engineering


- To increase the performance of your model, you can also experiment 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 

with X-axis as `n_estimators`, Y-axis as `max_depth`, and Z-axis as `AUC Score`, we have given the notebook which explains how to plot this 3d plot, you can find it in the same drive [3d_scatter_plot.ipynb](#)

(or)

- 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 

[seaborn heat maps](#) with rows as `n_estimators`, columns as `max_depth`, and values inside the cell representing **AUC Score**
- You choose either of the plotting techniques out of 3d plot or heat map
- 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-leakage, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method `fit_transform()` on your train data, and apply the method `transform()` on cv/test data.
4. For more details please go through this [link](#).

In [13]:

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import roc_auc_score
import scikitplot as skplt
from cycler import cycler
from sklearn.model_selection import GridSearchCV
from sklearn import svm
from sklearn.metrics import
accuracy_score, precision_score, recall_score, confusion_matrix, classification_report, f1_score
%matplotlib inline
import warnings
```

```
warnings.filterwarnings("ignore")
from wordcloud import WordCloud
import xgboost as xgb
```

[5.1] Applying RF

[5.1.1] Applying Random Forests on BOW, SET 1

In [14]:

```
X_train = X_train_BOW
X_test = X_test_BOW
```

In [15]:

```
max_depths = [2,4,6,9,11]
base_learners = [1, 5, 10, 50, 100]
param_grid = {'max_depth': max_depths, 'n_estimators': base_learners}

model = GridSearchCV(RandomForestClassifier(class_weight='balanced'), param_grid, scoring =
'roc_auc', cv=3, n_jobs = -1)
model.fit(X_train, y_train)
print("Model with best parameters :\n", model.best_estimator_)
print("Accuracy of the model : ", model.score(X_train, y_train))
```

Model with best parameters :

```
RandomForestClassifier(bootstrap=True, class_weight='balanced',
                        criterion='gini', max_depth=11, max_features='auto',
                        max_leaf_nodes=None, min_impurity_decrease=0.0,
                        min_impurity_split=None, min_samples_leaf=1,
                        min_samples_split=2, min_weight_fraction_leaf=0.0,
                        n_estimators=100, n_jobs=None, oob_score=False,
                        random_state=None, verbose=0, warm_start=False)
```

Accuracy of the model : 0.9219305673697699

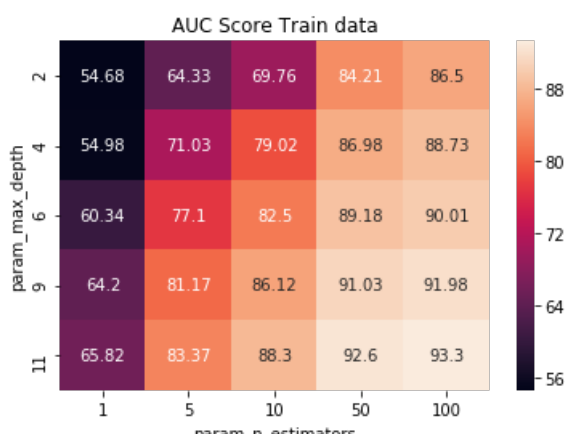
In [16]:

```
dataframe = pd.DataFrame(model.cv_results_) # model.cv_results_ : gives the results after fitting
the model
#Storing it into the dataframe and later plotting it into heatmap
```

In [17]:

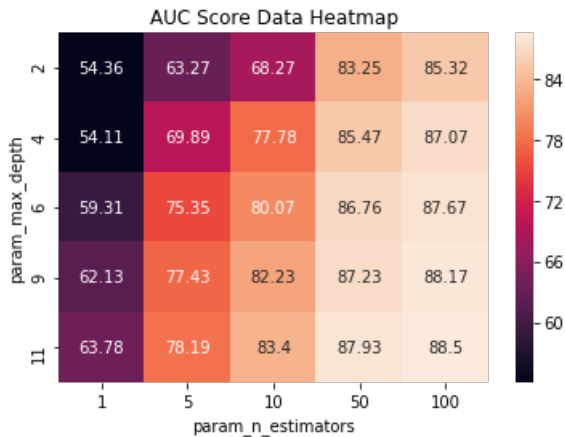
```
# Train Data AUC Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
                               'param_n_estimators']).max()

max_scores = max_scores.unstack()[['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_train_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score Train data')
plt.show()
```



In [18]:

```
# CV Data AUC Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
                                'param_n_estimators']).max()
max_scores = max_scores.unstack()[['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_test_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score Data Heatmap')
plt.show()
```



In [19]:

```
optimal_depth = 11
optimal_estimators = 100
```

In [20]:

```
lr = RandomForestClassifier(n_estimators=optimal_estimators, max_depth=optimal_depth, class_weight
                             ='balanced')
lr.fit(X_train_BOW,y_train)
pred = lr.predict(X_test_BOW)

print("***Test Data Report***")
print("Best max_depth = ",optimal_depth)
print("Best Base Learners = ",optimal_estimators)
fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[: ,1],pos_label="positive")
auc = metrics.auc(fpr, tpr)
print("AUC = ",auc*100)
skplt.metrics.plot_confusion_matrix(y_test, pred)
plt.show()

fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[: ,1],pos_label="positive")
fpr2, tpr2, threshold2 = metrics.roc_curve(y_train, lr.predict_proba(X_train)
[: ,1],pos_label="positive")

roc_auc = metrics.auc(fpr, tpr)
roc_auc2 = metrics.auc(fpr2, tpr2)

# method 1: plt
import matplotlib.pyplot as plt
f, ax = plt.subplots()
plt.title('Receiver Operating Characteristic Curve')
cy = cycler('color', ['red', 'green', 'blue'])
ax.set_prop_cycle(cy)
ax.plot(fpr, tpr, label = 'AUC = %0.2f' % roc_auc)
ax.plot(fpr2, tpr2, label = 'AUC = %0.2f' % roc_auc2)
plt.legend(['TEST Data', 'TRAIN Data'],loc = 'lower right')

ax.plot([0, 1], [0, 1])
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.xlabel('False Positive Rate')
```

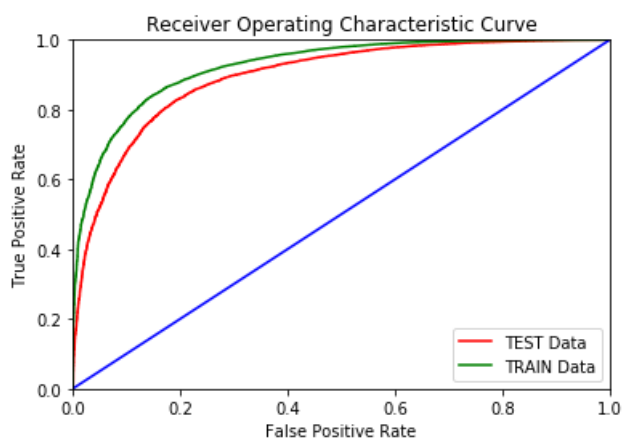
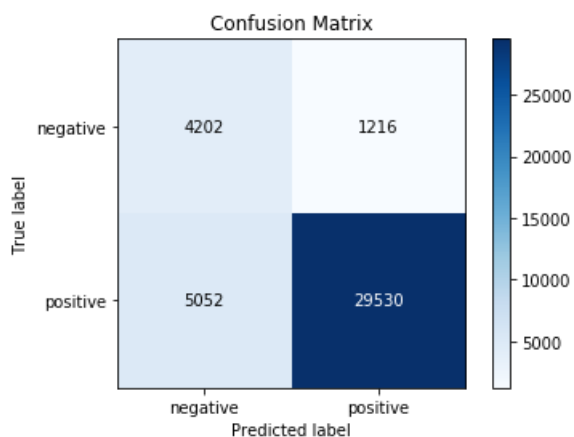
```
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()
```

Test Data Report

Best max_depth = 11

Best Base Learners = 100

AUC = 89.24354638689829



[5.1.2] Wordcloud of top 20 important features from SET 1

In [21]:

```
feature_name = count_vect.get_feature_names()
w = lr.feature_importances_
weight=w.reshape(-1)
sorted_feature = np.argsort(weight)
top_20_positive_feature=sorted_feature[:-20:-1]
```

In [22]:

```
top_feature_names = []
print("Top 20 features :")
print("-----")
for i in top_20_positive_feature:
    print("%s\t-->\t%f"%(feature_name[i],weight[i]))
    top_feature_names.append(feature_name[i])
```

Top 20 features :

great --> 0.063241

best --> 0.039361

disappoint --> 0.036742

love --> 0.030344

would --> 0.026406

perfect --> 0.026055

bad --> 0.022802

```

delici --> 0.022714
favorit --> 0.019374
thought --> 0.018315
money --> 0.017909
terribl --> 0.016676
worst --> 0.013297
horribl --> 0.012668
nice --> 0.012266
return --> 0.012128
easi --> 0.011857
wast --> 0.011355
aw --> 0.010347

```

In [23]:

```

#convert list to string and generate
unique_string=(" ").join(top_feature_names)
wordcloud = WordCloud(width = 1000, height = 400,background_color ='black').generate(unique_string)
plt.figure(figsize = (10, 8), facecolor = None)
plt.imshow(wordcloud)
plt.axis("off")
plt.savefig("your_file_name"+" .png", bbox_inches='tight')
plt.show()
plt.close()

```



[5.1.3] Applying Random Forests on TFIDF, SET 2

In [24]:

```

X_train = X_train_tfidf
X_test = X_test_tfidf

```

In [25]:

```

max_depths = [2,4,6,9,11]
base_learners = [1, 5, 10, 50, 100]
param_grid = {'max_depth': max_depths, 'n_estimators':base_learners}

model = GridSearchCV(RandomForestClassifier(class_weight='balanced'), param_grid, scoring =
'roc_auc', cv=3 , n_jobs = -1)
model.fit(X_train, y_train)
print("Model with best parameters :\n",model.best_estimator_)
print("Accuracy of the model : ",model.score(X_train, y_train))

```

Model with best parameters :

```

RandomForestClassifier(bootstrap=True, class_weight='balanced',
criterion='gini', max_depth=11, max_features='auto',
max_leaf_nodes=None, min_impurity_decrease=0.0,
min_impurity_split=None, min_samples_leaf=1,
min_samples_split=2, min_weight_fraction_leaf=0.0,
n_estimators=100, n_jobs=None, oob_score=False,
random_state=None, verbose=0, warm_start=False)

```

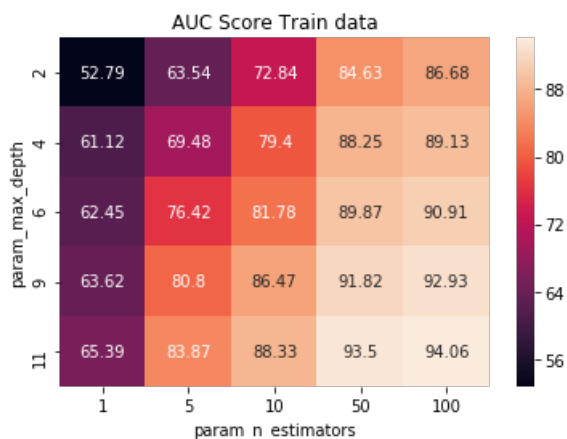
Accuracy of the model : 0.9305428376505988

In [26]:

```
dataframe = pd.DataFrame(model.cv_results_) # model.cv_results_ : gives the results after fitting the model
#Storing it into the dataframe and later plotting it into heatmap
```

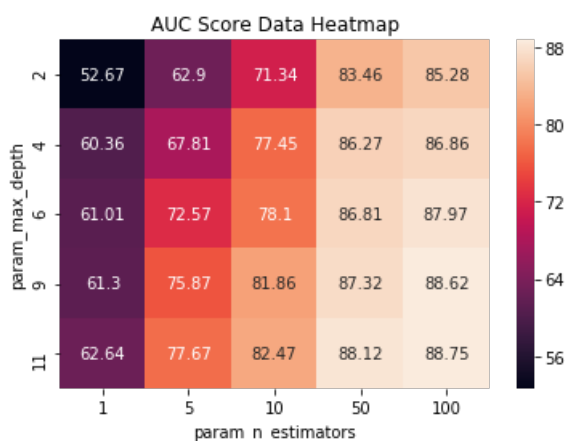
In [27]:

```
# Train Data AUC Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
                               'param_n_estimators']).max()
max_scores = max_scores.unstack()[['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_train_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score Train data')
plt.show()
```



In [28]:

```
# CV Data AUC Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
                               'param_n_estimators']).max()
max_scores = max_scores.unstack()[['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_test_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score Data Heatmap')
plt.show()
```



In [29]:

```
optimal_depth = 11
optimal_estimators = 100
```

In [30]:

```
lr = RandomForestClassifier(n_estimators=optimal_estimators, max_depth=optimal_depth, class_weight='balanced')
```

```

balanced ,
lr.fit(X_train,y_train)
pred = lr.predict(X_test)

print("***Test Data Report***")
print("Best max_depth = ",optimal_depth)
print("Best Base Learners = ",optimal_estimators)
fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[: ,1],pos_label="positive")
auc = metrics.auc(fpr, tpr)
print("AUC = ",auc*100)
skplt.metrics.plot_confusion_matrix(y_test, pred)
plt.show()

fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[: ,1],pos_label="positive")
fpr2, tpr2, threshold2 = metrics.roc_curve(y_train, lr.predict_proba(X_train)
[: ,1],pos_label="positive")

roc_auc = metrics.auc(fpr, tpr)
roc_auc2 = metrics.auc(fpr2, tpr2)

# method 1: plt
import matplotlib.pyplot as plt
f, ax = plt.subplots()
plt.title('Receiver Operating Characteristic Curve')
cy = cycler('color', ['red', 'green', 'blue'])
ax.set_prop_cycle(cy)
ax.plot(fpr, tpr, label = 'AUC = %0.2f' % roc_auc)
ax.plot(fpr2, tpr2, label = 'AUC = %0.2f' % roc_auc2)
plt.legend(['TEST Data', 'TRAIN Data'],loc = 'lower right')

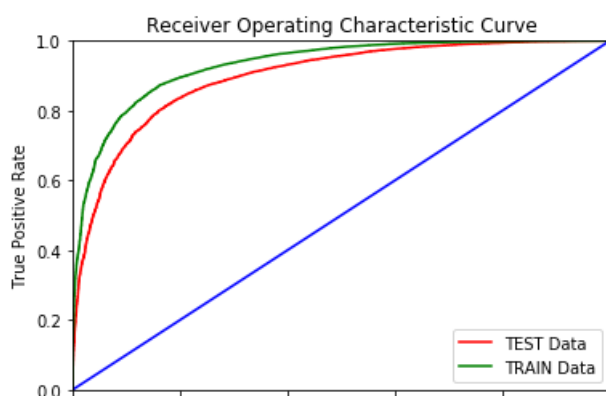
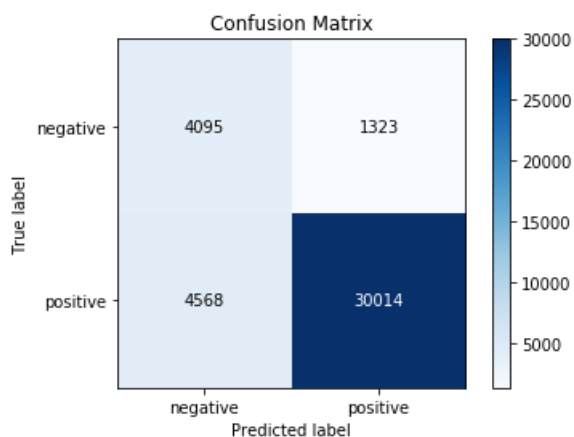
ax.plot([0, 1], [0, 1])
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()

```

```

***Test Data Report***
Best max_depth = 11
Best Base Learners = 100
AUC = 89.62925766458456

```



0.0 0.2 0.4 0.6 0.8 1.0
False Positive Rate

[5.1.4] Wordcloud of top 20 important features from SET 2

In [31]:

```
feature_name = tf_idf_vect.get_feature_names()
w = lr.feature_importances_
weight=w.reshape(-1)
sorted_feature = np.argsort(weight)
top_20_positive_feature=sorted_feature[:-20:-1]
```

In [32]:

```
top_feature_names = []
print("Top 20 features :")
print("-----")
for i in top_20_positive_feature:
    print("%s\t-->\t%f"%(feature_name[i],weight[i]))
    top_feature_names.append(feature_name[i])
```

```
Top 20 features :
-----
great --> 0.061601
disappoint --> 0.049728
delici --> 0.034036
love --> 0.033465
best --> 0.030766
money --> 0.024871
perfect --> 0.020179
would --> 0.018817
didnt --> 0.018099
bad --> 0.016966
easi --> 0.016554
return --> 0.015728
worst --> 0.015012
favorit --> 0.014654
aw --> 0.013993
horribl --> 0.013772
receiv --> 0.013229
mayb --> 0.013184
find --> 0.012421
```

In [33]:

```
#convert list to string and generate
unique_string=(" ").join(top_feature_names)
wordcloud = WordCloud(width = 1000, height = 400,background_color ='black').generate(unique_string)
plt.figure(figsize = (10, 8), facecolor = None)
plt.imshow(wordcloud)
plt.axis("off")
plt.savefig("your_file_name"+" .png", bbox_inches='tight')
plt.show()
plt.close()
```



[5.1.5] Applying Random Forests on AVG W2V, SET 3

In [63]:

```
X_train = sent_vectors_train
X_test = sent_vectors_test
```

In [64]:

```
max_depths = [2,4,6,9,11]
base_learners = [1, 5, 10, 50, 100]
param_grid = {'max_depth': max_depths, 'n_estimators':base_learners}

model = GridSearchCV(RandomForestClassifier(class_weight='balanced'), param_grid, scoring =
'roc_auc', cv=3, n_jobs = -1)
model.fit(X_train, y_train)
print("Model with best parameters :\n",model.best_estimator_)
print("Accuracy of the model : ",model.score(X_train, y_train))
```

Model with best parameters :

```
RandomForestClassifier(bootstrap=True, class_weight='balanced',
criterion='gini', max_depth=11, max_features='auto',
max_leaf_nodes=None, min_impurity_decrease=0.0,
min_impurity_split=None, min_samples_leaf=1,
min_samples_split=2, min_weight_fraction_leaf=0.0,
n_estimators=100, n_jobs=None, oob_score=False,
random_state=None, verbose=0, warm_start=False)
```

Accuracy of the model : 0.9713874175970296

In [65]:

```
dataframe = pd.DataFrame(model.cv_results_) # model.cv_results_ : gives the results after fitting
the model
#Storing it into the dataframe and later plotting it into heatmap
```

In [66]:

```
# Train Data AUC Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
'param_n_estimators']).max()
max_scores = max_scores.unstack()[['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_train_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score Train data')
plt.show()
```



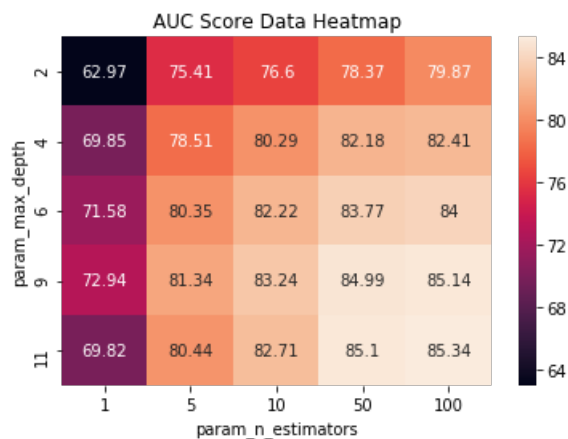
In [67]:

```
# CV Data AUC Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
```

```

max_scores = max_scores.unstack()['param_n_estimators'].max()
sns.heatmap(max_scores.mean_test_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score Data Heatmap')
plt.show()

```



In [68]:

```

optimal_depth = 11
optimal_estimators = 100

```

In [69]:

```

lr = RandomForestClassifier(n_estimators=optimal_estimators, max_depth=optimal_depth, class_weight
='balanced')
lr.fit(X_train,y_train)
pred = lr.predict(X_test)

print("***Test Data Report***")
print("Best max_depth = ",optimal_depth)
print("Best Base Learners = ",optimal_estimators)
fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[: ,1],pos_label="positive")
auc = metrics.auc(fpr, tpr)
print("AUC = ",auc*100)
skplt.metrics.plot_confusion_matrix(y_test, pred)
plt.show()

fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[: ,1],pos_label="positive")
fpr2, tpr2, threshold2 = metrics.roc_curve(y_train, lr.predict_proba(X_train)
[: ,1],pos_label="positive")

roc_auc = metrics.auc(fpr, tpr)
roc_auc2 = metrics.auc(fpr2, tpr2)

# method I: plt
import matplotlib.pyplot as plt
f, ax = plt.subplots()
plt.title('Receiver Operating Characteristic Curve')
cy = cycler('color', ['red', 'green', 'blue'])
ax.set_prop_cycle(cy)
ax.plot(fpr, tpr, label = 'AUC = %0.2f' % roc_auc)
ax.plot(fpr2, tpr2, label = 'AUC = %0.2f' % roc_auc2)
plt.legend(['TEST Data', 'TRAIN Data'],loc = 'lower right')

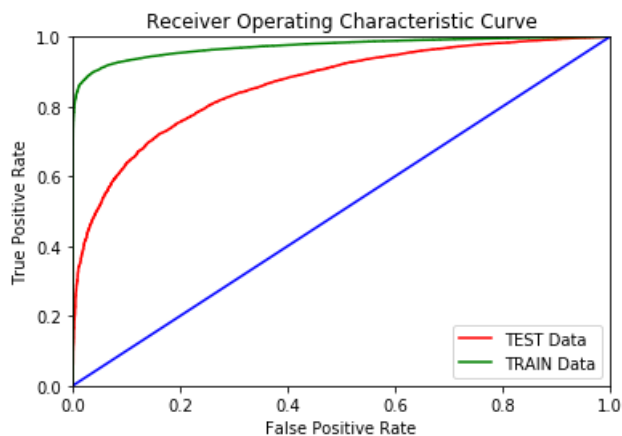
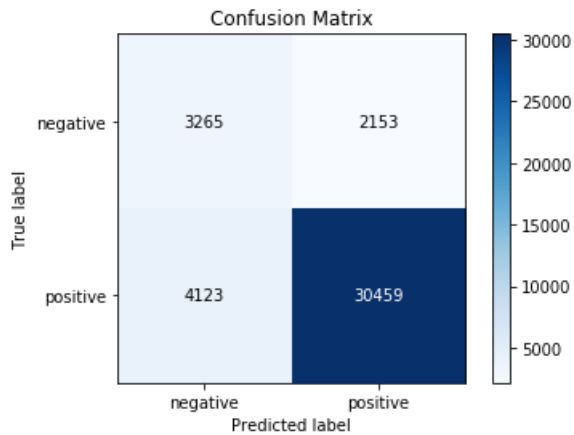
ax.plot([0, 1], [0, 1])
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()

```

```

***Test Data Report***
Best max_depth = 11
Best Base Learners = 100
AUC = 86.06265335952645

```



[5.1.6] Applying Random Forests on TFIDF W2V, SET 4

In [71]:

```
X_train = tfidf_sent_vectors_train
X_test = tfidf_sent_vectors_test
```

In [72]:

```
max_depths = [2,4,6,9,11]
base_learners = [1, 5, 10, 50, 100]
param_grid = {'max_depth': max_depths, 'n_estimators': base_learners}

model = GridSearchCV(RandomForestClassifier(class_weight='balanced'), param_grid, scoring =
'roc_auc', cv=3, n_jobs = -1)
model.fit(X_train, y_train)
print("Model with best parameters :\n", model.best_estimator_)
print("Accuracy of the model : ", model.score(X_train, y_train))
```

Model with best parameters :

```
RandomForestClassifier(bootstrap=True, class_weight='balanced',
criterion='gini', max_depth=11, max_features='auto',
max_leaf_nodes=None, min_impurity_decrease=0.0,
min_impurity_split=None, min_samples_leaf=1,
min_samples_split=2, min_weight_fraction_leaf=0.0,
n_estimators=100, n_jobs=None, oob_score=False,
random_state=None, verbose=0, warm_start=False)
```

Accuracy of the model : 0.9627524835915471

In [73]:

```
dataframe = pd.DataFrame(model.cv_results_) # model.cv_results_ : gives the results after fitting
the model
#Storing it into the dataframe and later plotting it into heatmap
```

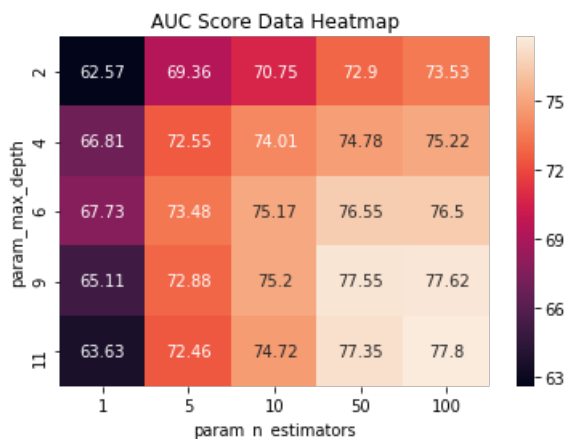
In [74]:

```
# Train Data AUC Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
                               'param_n_estimators']).max()
max_scores = max_scores.unstack()[['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_train_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score Train data')
plt.show()
```



In [75]:

```
# CV Data AUC Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
                               'param_n_estimators']).max()
max_scores = max_scores.unstack()[['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_test_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score Data Heatmap')
plt.show()
```



In [76]:

```
optimal_depth = 11
optimal_estimators = 100
```

In [77]:

```
lr = RandomForestClassifier(n_estimators=optimal_estimators, max_depth=optimal_depth, class_weight='balanced')
lr.fit(X_train,y_train)
pred = lr.predict(X_test)

print("***Test Data Report***")
print("Best max_depth = ",optimal_depth)
print("Best Base Learners = " ,optimal_estimators)
```

```

print('Best Base Learners = ', optimal_learners)
fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[: ,1], pos_label="positive")
auc = metrics.auc(fpr, tpr)
print("AUC = ", auc*100)
skplt.metrics.plot_confusion_matrix(y_test, pred)
plt.show()

fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[: ,1], pos_label="positive")
fpr2, tpr2, threshold2 = metrics.roc_curve(y_train, lr.predict_proba(X_train)
[: ,1], pos_label="positive")

roc_auc = metrics.auc(fpr, tpr)
roc_auc2 = metrics.auc(fpr2, tpr2)

# method I: plt
import matplotlib.pyplot as plt
f, ax = plt.subplots()
plt.title('Receiver Operating Characteristic Curve')
cy = cycler('color', ['red', 'green', 'blue'])
ax.set_prop_cycle(cy)
ax.plot(fpr, tpr, label = 'AUC = %0.2f' % roc_auc)
ax.plot(fpr2, tpr2, label = 'AUC = %0.2f' % roc_auc2)
plt.legend(['TEST Data', 'TRAIN Data'], loc = 'lower right')

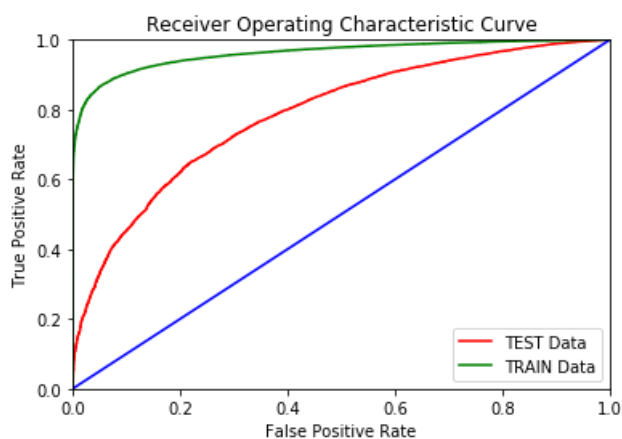
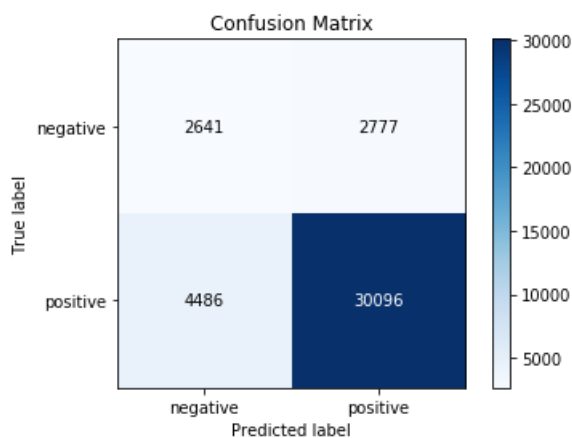
ax.plot([0, 1], [0, 1])
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()

```

```

***Test Data Report***
Best max_depth = 11
Best Base Learners = 100
AUC = 78.63989963647266

```



[5.2] Applying GBDT using XGBOOST

[5.2.1] Applying XGBOOST on BOW, SET 1

In [35]:

```
X_train = X_train_BOW
X_test = X_test_BOW
```

In [36]:

```
max_depths = [2,4,6,9,11]
base_learners = [1, 5, 10, 50, 100]
param_grid = {'max_depth': max_depths, 'n_estimators':base_learners}

# scale_pos_weight=1 for Balancing of positive and negative weights
# https://xgboost.readthedocs.io/en/latest/python/python_api.html
model = GridSearchCV(xgb.XGBClassifier(scale_pos_weight=1), param_grid, scoring = 'roc_auc', cv=3 ,
n_jobs = -1)
model.fit(X_train, y_train)
print("Model with best parameters :\n",model.best_estimator_)
print("Accuracy of the model : ",model.score(X_train, y_train))
```

Model with best parameters :

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
max_depth=11, min_child_weight=1, missing=None, n_estimators=100,
n_jobs=1, nthread=None, objective='binary:logistic', random_state=0,
reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
silent=True, subsample=1)
```

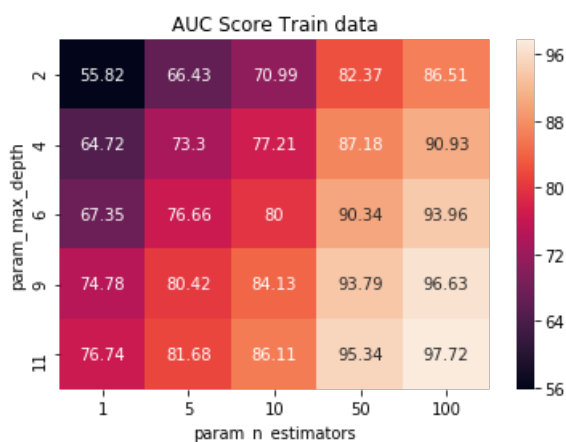
Accuracy of the model : 0.9715711023560167

In [37]:

```
dataframe = pd.DataFrame(model.cv_results_) # model.cv_results_ : gives the results after fitting
the model
#Storing it into the dataframe and later plotting it into heatmap
```

In [38]:

```
# Train Data Auc Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
                              'param_n_estimators']).max()
max_scores = max_scores.unstack() [['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_train_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score Train data')
plt.show()
```



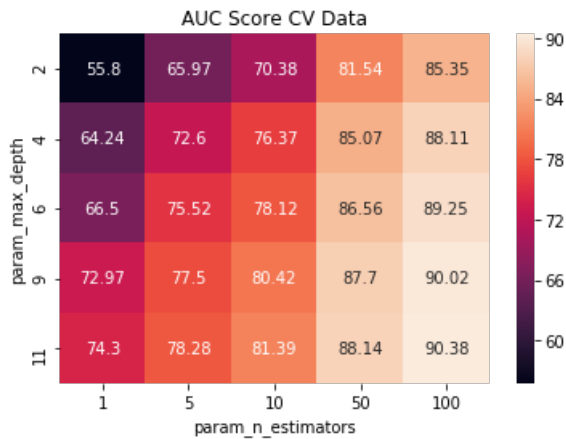
In [39]:

```
# CV Data Auc Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
                              'param_n_estimators']).max()
max_scores = max_scores.unstack() [['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_test_score*100, annot=True, fmt='.4g');
```

```

xgb.XGBClassifier(max_depth=optimal_depth, n_estimators=optimal_estimators, scale_pos_weight=1)
ax = plt.axes()
ax.set_title('AUC Score CV Data')
plt.show()

```



In [40]:

```

optimal_depth = 11
optimal_estimators = 100

```

In [41]:

```

lr = xgb.XGBClassifier(max_depth=optimal_depth, n_estimators=optimal_estimators, scale_pos_weight=1)
lr.fit(X_train,y_train)
pred = lr.predict(X_test)

print("***Test Data Report***")
print("Best max_depth = ",optimal_depth)
print("Best Base Learners = ",optimal_estimators)
fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[: ,1],pos_label="positive")
auc = metrics.auc(fpr, tpr)
print("AUC = ",auc*100)
skplt.metrics.plot_confusion_matrix(y_test, pred)
plt.show()

fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[: ,1],pos_label="positive")
fpr2, tpr2, threshold2 = metrics.roc_curve(y_train, lr.predict_proba(X_train)[: ,1],pos_label="positive")

roc_auc = metrics.auc(fpr, tpr)
roc_auc2 = metrics.auc(fpr2, tpr2)

# method 1: plt
import matplotlib.pyplot as plt
f, ax = plt.subplots()
plt.title('Receiver Operating Characteristic Curve')
cy = cycler('color', ['red', 'green', 'blue'])
ax.set_prop_cycle(cy)
ax.plot(fpr, tpr, label = 'AUC = %0.2f' % roc_auc)
ax.plot(fpr2, tpr2, label = 'AUC = %0.2f' % roc_auc2)
plt.legend(['TEST Data', 'TRAIN Data'],loc = 'lower right')

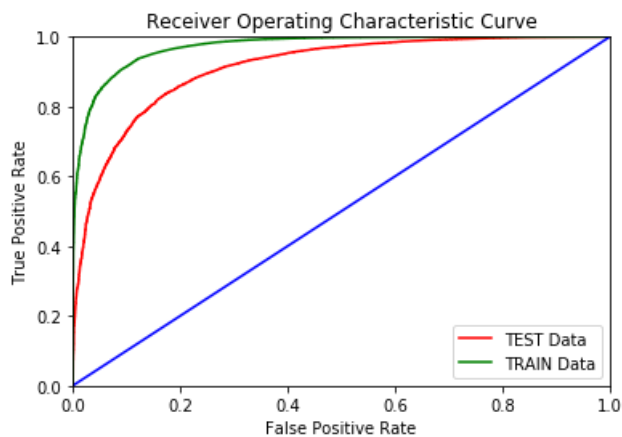
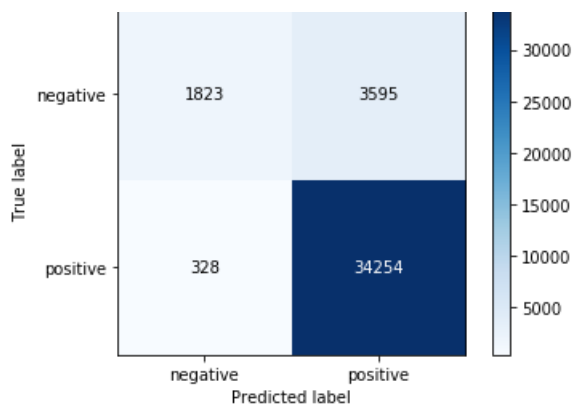
ax.plot([0, 1], [0, 1])
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()

```

```

***Test Data Report***
Best max_depth = 11
Best Base Learners = 100
AUC = 91.06288536622976

```



[5.2.2] Applying XGBOOST on TFIDF, SET 2

In [43]:

```
X_train = X_train_tfidf
X_test = X_test_tfidf
```

In [44]:

```
max_depths = [2,4,6,9,11]
base_learners = [1, 5, 10, 50, 100]
param_grid = {'max_depth': max_depths, 'n_estimators': base_learners}

# scale_pos_weight=1 for Balancing of positive and negative weights
# https://xgboost.readthedocs.io/en/latest/python/python_api.html
model = GridSearchCV(xgb.XGBClassifier(scale_pos_weight=1), param_grid, scoring = 'roc_auc', cv=3 ,
n_jobs = -1)
model.fit(X_train, y_train)
print("Model with best parameters :\n", model.best_estimator_)
print("Accuracy of the model : ", model.score(X_train, y_train))
```

Model with best parameters :

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
max_depth=11, min_child_weight=1, missing=None, n_estimators=100,
n_jobs=1, nthread=None, objective='binary:logistic', random_state=0,
reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
silent=True, subsample=1)
```

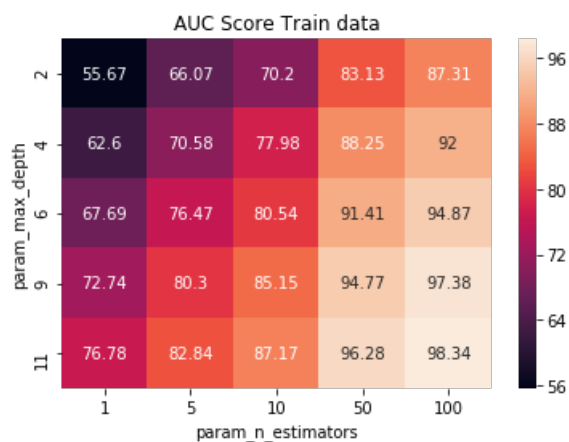
Accuracy of the model : 0.9783129818009755

In [45]:

```
dataframe = pd.DataFrame(model.cv_results_) # model.cv_results_ : gives the results after fitting
the model
#Storing it into the dataframe and later plotting it into heatmap
```

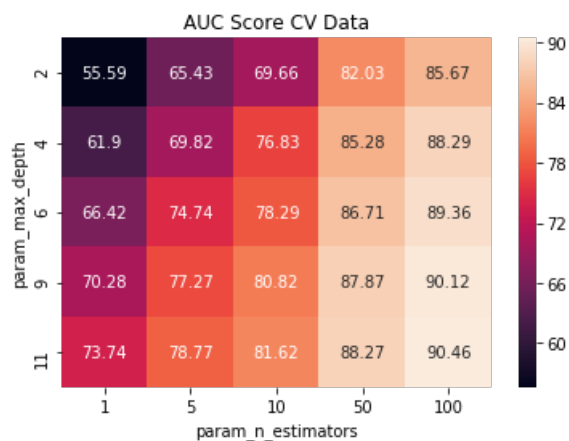
In [46]:

```
# Train Data AUC Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
                                'param_n_estimators']).max()
max_scores = max_scores.unstack()[['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_train_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score Train data')
plt.show()
```



In [47]:

```
# CV Data AUC Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
                                'param_n_estimators']).max()
max_scores = max_scores.unstack()[['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_test_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score CV Data')
plt.show()
```



In [48]:

```
optimal_depth = 11
optimal_estimators = 100
```

In [49]:

```
lr = xgb.XGBClassifier(max_depth=optimal_depth, n_estimators=optimal_estimators, scale_pos_weight=1)
lr.fit(X_train,y_train)
pred = lr.predict(X_test)

print("***Test Data Report***")
print("Best max_depth = ",optimal_depth)
print("Best Base Learners = ",optimal_estimators)
fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[:,1],pos_label="positive")
auc = metrics.auc(fpr, tpr)
print("AUC = " + auc*100)
```

```

print(AUC - ,auc_100)
skplt.metrics.plot_confusion_matrix(y_test, pred)
plt.show()

fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[: ,1],pos_label="positive")
fpr2, tpr2, threshold2 = metrics.roc_curve(y_train, lr.predict_proba(X_train)
[: ,1],pos_label="positive")

roc_auc = metrics.auc(fpr, tpr)
roc_auc2 = metrics.auc(fpr2, tpr2)

# method 1: plt
import matplotlib.pyplot as plt
f, ax = plt.subplots()
plt.title('Receiver Operating Characteristic Curve')
cy = cycler('color', ['red', 'green', 'blue'])
ax.set_prop_cycle(cy)
ax.plot(fpr, tpr, label = 'AUC = %0.2f' % roc_auc)
ax.plot(fpr2, tpr2, label = 'AUC = %0.2f' % roc_auc2)
plt.legend(['TEST Data', 'TRAIN Data'],loc = 'lower right')

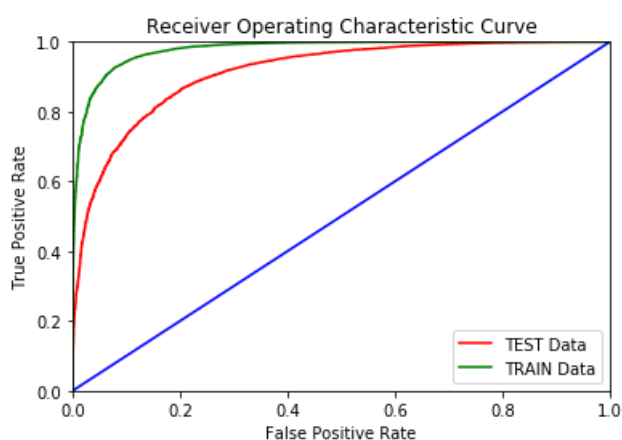
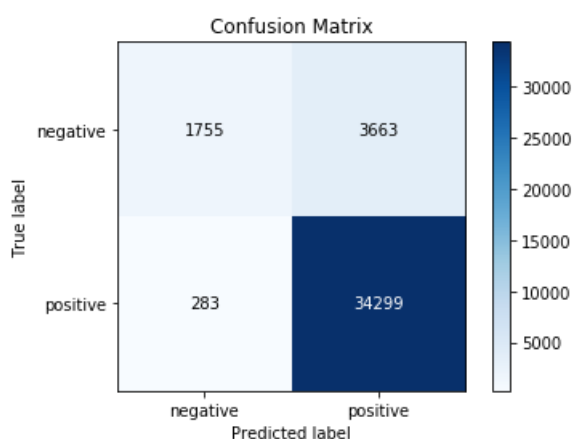
ax.plot([0, 1], [0, 1])
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()

```

```

***Test Data Report***
Best max_depth = 11
Best Base Learners = 100
AUC = 91.27022768082172

```



[5.2.3] Applying XGBOOST on AVG W2V, SET 3

In [78]:

```

X_train = sent_vectors_train
X_test = sent_vectors_test

```

In [80]:

```
max_depths = [2,4,6,9,11]
base_learners = [1, 5, 10, 50, 100]
param_grid = {'max_depth': max_depths, 'n_estimators': base_learners}

# scale_pos_weight=1 for Balancing of positive and negative weights
# https://xgboost.readthedocs.io/en/latest/python/python_api.html
model = GridSearchCV(xgb.XGBClassifier(scale_pos_weight=1), param_grid, scoring = 'roc_auc', cv=3,
n_jobs = -1)
model.fit(np.array(X_train), y_train)
print("Model with best parameters :\n", model.best_estimator_)
```

Model with best parameters :

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
max_depth=6, min_child_weight=1, missing=None, n_estimators=100,
n_jobs=1, nthread=None, objective='binary:logistic', random_state=0,
reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
silent=True, subsample=1)
```

In [81]:

```
dataframe = pd.DataFrame(model.cv_results_) # model.cv_results_ : gives the results after fitting
the model
# Storing it into the dataframe and later plotting it into heatmap
```

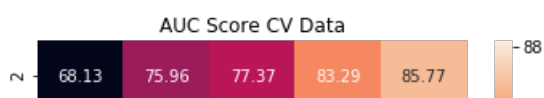
In [82]:

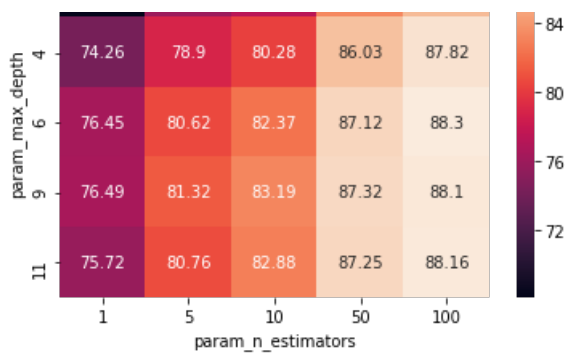
```
# Train Data AUC Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
                               'param_n_estimators']).max()
max_scores = max_scores.unstack() [['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_train_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score Train data')
plt.show()
```



In [83]:

```
# CV Data AUC Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
                               'param_n_estimators']).max()
max_scores = max_scores.unstack() [['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_test_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score CV Data')
plt.show()
```





In [84]:

```
optimal_depth = 6
optimal_estimators = 100
```

In [85]:

```
lr = xgb.XGBClassifier(max_depth=optimal_depth, n_estimators=optimal_estimators, scale_pos_weight=1)
lr.fit(np.array(X_train), y_train)
pred = lr.predict(np.array(X_test))

print("***Test Data Report***")
print("Best max_depth = ", optimal_depth)
print("Best Base Learners = ", optimal_estimators)
fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[:,1], pos_label="positive")
auc = metrics.auc(fpr, tpr)
print("AUC = ", auc*100)
skplt.metrics.plot_confusion_matrix(y_test, pred)
plt.show()

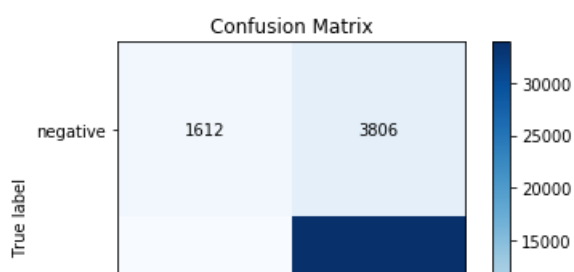
fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[:,1], pos_label="positive")
fpr2, tpr2, threshold2 = metrics.roc_curve(y_train, lr.predict_proba(X_train)[:,1], pos_label="positive")

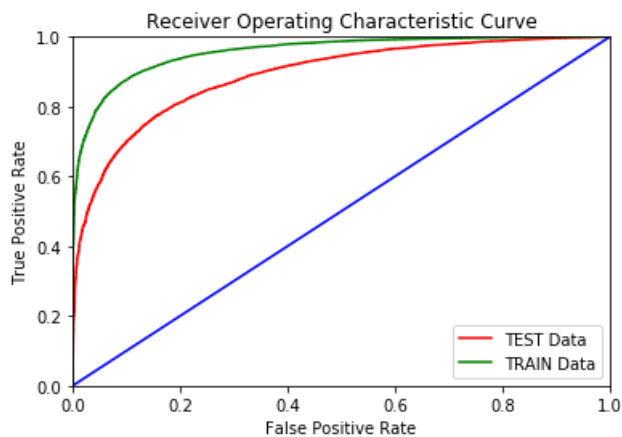
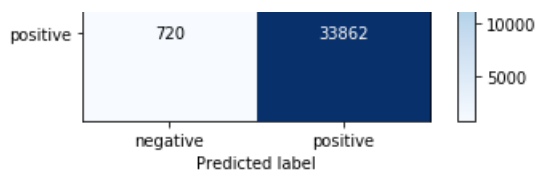
roc_auc = metrics.auc(fpr, tpr)
roc_auc2 = metrics.auc(fpr2, tpr2)

# method 1: plt
import matplotlib.pyplot as plt
f, ax = plt.subplots()
plt.title('Receiver Operating Characteristic Curve')
cy = cycycler('color', ['red', 'green', 'blue'])
ax.set_prop_cycle(cy)
ax.plot(fpr, tpr, label = 'AUC = %0.2f' % roc_auc)
ax.plot(fpr2, tpr2, label = 'AUC = %0.2f' % roc_auc2)
plt.legend(['TEST Data', 'TRAIN Data'], loc = 'lower right')

ax.plot([0, 1], [0, 1])
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()
```

```
***Test Data Report***
Best max_depth = 6
Best Base Learners = 100
AUC = 88.91754654688523
```





[5.2.4] Applying XGBOOST on TFIDF W2V, SET 4

In [86]:

```
X_train = tfidf_sent_vectors_train
X_test = tfidf_sent_vectors_test
```

In [87]:

```
max_depths = [2,4,6,9,11]
base_learners = [1, 5, 10, 50, 100]
param_grid = {'max_depth': max_depths, 'n_estimators': base_learners}

# scale_pos_weight=1 for Balancing of positive and negative weights
# https://xgboost.readthedocs.io/en/latest/python/python_api.html
model = GridSearchCV(xgb.XGBClassifier(scale_pos_weight=1), param_grid, scoring = 'roc_auc', cv=3 ,
n_jobs = -1)
model.fit(np.array(X_train), y_train)
print("Model with best parameters :\n", model.best_estimator_)
print("Accuracy of the model : ", model.score(X_train, y_train))
```

Model with best parameters :

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
max_depth=6, min_child_weight=1, missing=None, n_estimators=100,
n_jobs=1, nthread=None, objective='binary:logistic', random_state=0,
reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
silent=True, subsample=1)
```

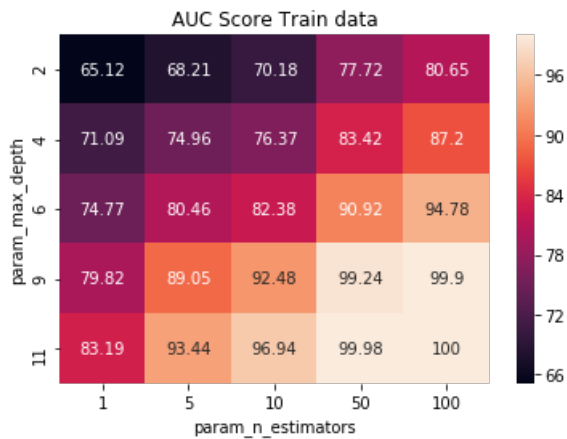
Accuracy of the model : 0.9258647285865835

In [88]:

```
dataframe = pd.DataFrame(model.cv_results_) # model.cv_results_ : gives the results after fitting
the model
#Storing it into the dataframe and later plotting it into heatmap
```

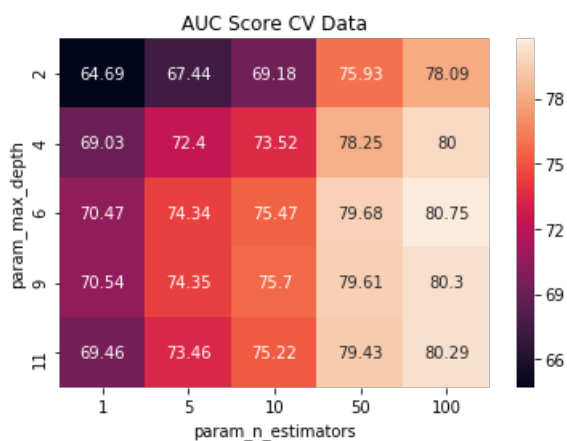
In [89]:

```
# Train Data Auc Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
'param_n_estimators']).max()
max_scores = max_scores.unstack()[['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_train_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score Train data')
plt.show()
```

In [90]:

```
# CV Data AUC Score Vs hyperparameter Heatmap
max_scores = dataframe.groupby(['param_max_depth',
                               'param_n_estimators']).max()
max_scores = max_scores.unstack()[['mean_test_score', 'mean_train_score']]
sns.heatmap(max_scores.mean_test_score*100, annot=True, fmt='.4g');
ax = plt.axes()
ax.set_title('AUC Score CV Data')
plt.show()
```



In [91]:

```
optimal_depth = 6
optimal_estimators = 100
```

In [93]:

```
lr = xgb.XGBClassifier(max_depth=optimal_depth, n_estimators=optimal_estimators, scale_pos_weight=1)
lr.fit(np.array(X_train), y_train)
pred = lr.predict(np.array(X_test))

print("***Test Data Report***")
print("Best max_depth = ", optimal_depth)
print("Best Base Learners = ", optimal_estimators)
fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[:,1], pos_label="positive")
auc = metrics.auc(fpr, tpr)
print("AUC = ", auc*100)
skplt.metrics.plot_confusion_matrix(y_test, pred)
plt.show()

fpr, tpr, threshold = metrics.roc_curve(y_test, lr.predict_proba(X_test)[:,1], pos_label="positive")
fpr2, tpr2, threshold2 = metrics.roc_curve(y_train, lr.predict_proba(X_train)[:,1], pos_label="positive")

roc_auc = metrics.auc(fpr, tpr)
```

```

roc_auc2 = metrics.auc(fpr2, tpr2)

# method 1: plt
import matplotlib.pyplot as plt
f, ax = plt.subplots()
plt.title('Receiver Operating Characteristic Curve')
cy = cycler('color', ['red', 'green', 'blue'])
ax.set_prop_cycle(cy)
ax.plot(fpr, tpr, label = 'AUC = %0.2f' % roc_auc)
ax.plot(fpr2, tpr2, label = 'AUC = %0.2f' % roc_auc2)
plt.legend(['TEST Data', 'TRAIN Data'],loc = 'lower right')

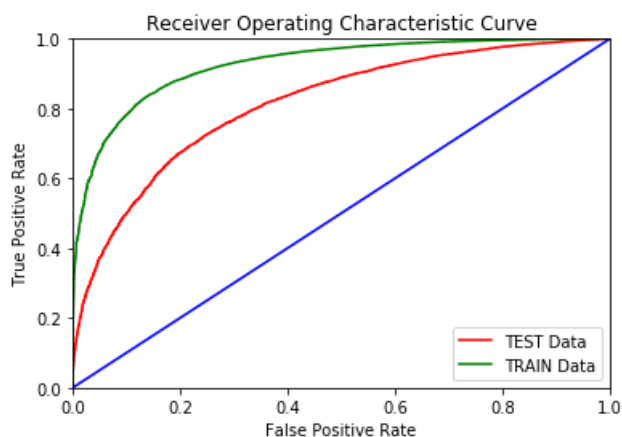
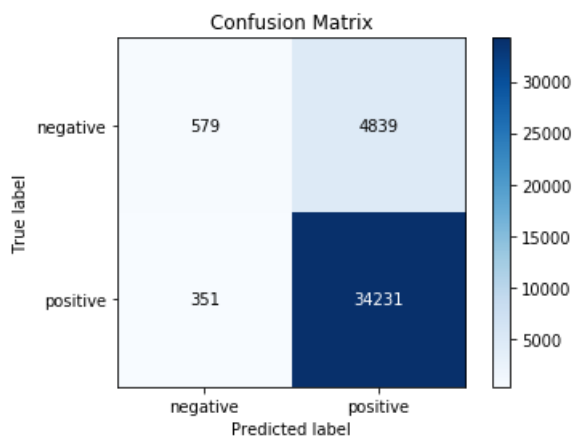
ax.plot([0, 1], [0, 1])
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()

```

```

***Test Data Report***
Best max_depth = 6
Best Base Learners = 100
AUC = 81.30034697571176

```



[6] Conclusions

In [97]:

```

#importing library
from prettytable import PrettyTable
x = PrettyTable()

#adding Field names
x.field_names = ["SL No.", "Vectorizer", "Model", "Best (max_depth)", "Best (base_learners)", "AUC"]

# adding row to table
x.title = 'Random Forest'
x.add_row(["1", "BOW", 'RandomForestClassifier', 11, 100, 89.2435])

```

```

x.add_row(["2","TFIDF",'RandomForestClassifier',11,100,89.6292])
x.add_row(["3","Avg-W2vec",'RandomForestClassifier',11,100,91.0628])
x.add_row(["4","TFIDF-W2vec",'RandomForestClassifier',11,100,78.6398])
x.add_row(["5","BOW",'XGBoost',11,100,91.0628])
x.add_row(["6","TFIDF",'XGBoost',11,100,91.2702])
x.add_row(["7","Avg-W2vec",'XGBoost',6,100,88.9176])
x.add_row(["8","TFIDF-W2vec",'XGBoost',6,100,81.3003])

```

```

#printing the table
print(x)

```

| SL No. | Vectorizer | Model | Best (max_depth) | Best (base_learners) | AUC |
|--------|-------------|------------------------|------------------|----------------------|---------|
| 1 | BOW | RandomForestClassifier | 11 | 100 | 89.2435 |
| 2 | TFIDF | RandomForestClassifier | 11 | 100 | 89.6292 |
| 3 | Avg-W2vec | RandomForestClassifier | 11 | 100 | 91.0628 |
| 4 | TFIDF-W2vec | RandomForestClassifier | 11 | 100 | 78.6398 |
| 5 | BOW | XGBoost | 11 | 100 | 91.0628 |
| 6 | TFIDF | XGBoost | 11 | 100 | 91.2702 |
| 7 | Avg-W2vec | XGBoost | 6 | 100 | 88.9176 |
| 8 | TFIDF-W2vec | XGBoost | 6 | 100 | 81.3003 |