

Amazon Fine Food Reviews Analysis

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

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

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

Number of reviews: 568,454

Number of users: 256,059

Number of products: 74,258

Timespan: Oct 1999 - Oct 2012

Number of Attributes/Columns in data: 10

Attribute Information:

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

Objective:

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

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

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

[1]. Reading Data

[1.1] Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

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

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

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

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

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

Number of data points in our data (50000, 10)

Out[3]:

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

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

```
In [5]: 1 print(display.shape)
        2 display.head()
```

(80668, 7)

Out[5]:

	UserId	ProductId	ProfileName	Time	Score	Text	COUNT(*)
0	#oc-R115TNMSPFT9I7	B007Y59HVM	Breyton	1331510400	2	Overall its just OK when considering the price...	2
1	#oc-R11D9D7SHXIJB9	B005HG9ET0	Louis E. Emory "hoppy"	1342396800	5	My wife has recurring extreme muscle spasms, u...	3
2	#oc-R11DNU2NBKQ23Z	B007Y59HVM	Kim Cieszykowski	1348531200	1	This coffee is horrible and unfortunately not ...	2
3	#oc-R11O5J5ZVQE25C	B005HG9ET0	Penguin Chick	1346889600	5	This will be the bottle that you grab from the...	3
4	#oc-R12KPBODL2B5ZD	B007OSBE1U	Christopher P. Presta	1348617600	1	I didnt like this coffee. Instead of telling y...	2

```
In [6]: 1 display[display['UserId'] == 'AZY10LLTJ71NX']
```

Out[6]:

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

```
In [7]: 1 display['COUNT(*)'].sum()
```

Out[7]: 393063

[2] Exploratory Data Analysis

[2.1] Data Cleaning: Deduplication

It is observed (as shown in the table below) that the reviews data had many duplicate entries. Hence it was necessary to remove duplicates in order to get unbiased results for the analysis of the data. Following is an example:

```
In [8]: 1 display= pd.read_sql_query("""
2 SELECT *
3 FROM Reviews
4 WHERE Score != 3 AND UserId="AR5J8UI46CURR"
5 ORDER BY ProductID
6 """, con)
7 display.head()
```

```
Out[8]:
```

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenon
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	

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

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

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

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

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

```
In [9]: 1 #Sorting data according to ProductId in ascending order
2 sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True, in
```

```
In [10]: 1 #Deduplication of entries
          2 final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"}
          3 final.shape
```

Out[10]: (46072, 10)

```
In [11]: 1 #Checking to see how much % of data still remains
          2 (final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

Out[11]: 92.144

Observation:- It was also seen that in two rows given below the value of HelpfulnessNumerator is greater than HelpfulnessDenominator which is not practically possible hence these two rows too are removed from calculations

```
In [12]: 1 display= pd.read_sql_query("""
          2 SELECT *
          3 FROM Reviews
          4 WHERE Score != 3 AND Id=44737 OR Id=64422
          5 ORDER BY ProductID
          6 """, con)
          7
          8 display.head()
```

Out[12]:

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



```
In [13]: 1 final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [14]: 1 #Before starting the next phase of preprocessing Lets see the number of entries
          2 print(final.shape)
          3
          4 #How many positive and negative reviews are present in our dataset?
          5 final['Score'].value_counts()
```

(46071, 10)

```
Out[14]: 1    38479
          0     7592
          Name: Score, dtype: int64
```

[3] Preprocessing

[3.1]. Preprocessing Review Text

Now that we have finished deduplication our data requires some preprocessing before we go on further with analysis and making the prediction model.

Hence in the Preprocessing phase we do the following in the order below:-

1. Begin by removing the html tags
2. Remove any punctuations or limited set of special characters like , or . or # etc.
3. Check if the word is made up of english letters and is not alpha-numeric
4. Check to see if the length of the word is greater than 2 (as it was researched that there is no adjective in 2-letters)
5. Convert the word to lowercase
6. Remove Stopwords
7. Finally Snowball Stemming the word (it was observed to be better than Porter Stemming)

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

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

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

=====

this is yummy, easy and unusual. it makes a quick, delicious pie, crisp or cobbler. home made is better, but a heck of a lot more work. this is great to have on hand for last minute dessert needs where you really want to impress with your creativity in cooking! recommended.

=====

Great flavor, low in calories, high in nutrients, high in protein! Usually protein powders are high priced and high in calories, this one is a great bargain and tastes great, I highly recommend for the lady gym rats, probably not "macho" enough for guys since it is soy based...

=====

For those of you wanting a high-quality, yet affordable green tea, you should definitely give this one a try. Let me first start by saying that everyone is looking for something different for their ideal tea, and I will attempt to briefly highlight what makes this tea attractive to a wide range of tea drinkers (whether you are a beginner or long-time tea enthusiast). I have gone through over 12 boxes of this tea myself, and highly recommend it for the following reasons:

-Quality: First, this tea offers a smooth quality without any harsh or bitter after tones, which often turns people off from many green teas. I've found my ideal brewing time to be between 3-5 minutes, giving you a light but flavorful cup of tea. However, if you get distracted or forget about your tea and leave it brewing for 20+ minutes like I sometimes do, the quality of this tea is such that you still get a smooth but deeper flavor without the bad after taste. The leaves themselves are whole leaves (not powdered stems, branches, etc commonly found in other brands), and the high-quality nylon bags also include chunks of tropical fruit and other discernible ingredients. This isn't your standard cheap paper bag with a mix of unknown ingredients that have been ground down to a fine powder, leaving you to wonder what it is you are actually drinking.

-Taste: This tea offers notes of real pineapple and other hints of tropical fruits, yet isn't sweet or artificially flavored. You have the foundation of a high-quality young hyson green tea for those true "tea flavor" lovers, yet the subtle hints of fruit make this a truly unique tea that I believe most will enjoy. If you want it sweet, you can add sugar, splenda, etc but this really is not necessary as this tea offers an inherent warmth of flavor through its ingredients.

-Price: This tea offers an excellent prod

uct at an exceptional price (especially when purchased at the prices Amazon offers). Compared to other brands which I believe to be of similar quality (Mighty Leaf, Rishi, Two Leaves, etc.), Revolution offers a superior product at an outstanding price. I have been purchasing this through Amazon for less per box than I would be paying at my local grocery store for Lipton, etc.

Overall, this is a wonderful tea that is comparable, and even better than, other teas that are priced much higher. It offers a well-balanced cup of green tea that I believe many will enjoy. In terms of taste, quality, and price, I would argue you won't find a better combination that that offered by Revolution's Tropical Green Tea.

=====

In [16]:

```
1 # remove urls from text python: https://stackoverflow.com/a/40823105/4084039
2 sent_0 = re.sub(r"http\S+", "", sent_0)
3 sent_1000 = re.sub(r"http\S+", "", sent_1000)
4 sent_150 = re.sub(r"http\S+", "", sent_1500)
5 sent_4900 = re.sub(r"http\S+", "", sent_4900)
6
7 print(sent_0)
```

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

```
In [17]: 1 # https://stackoverflow.com/questions/16206380/python-beautifulsoup-how-to-re
2 from bs4 import BeautifulSoup
3
4 soup = BeautifulSoup(sent_0, 'lxml')
5 text = soup.get_text()
6 print(text)
7 print("=="*50)
8
9 soup = BeautifulSoup(sent_1000, 'lxml')
10 text = soup.get_text()
11 print(text)
12 print("=="*50)
13
14 soup = BeautifulSoup(sent_1500, 'lxml')
15 text = soup.get_text()
16 print(text)
17 print("=="*50)
18
19 soup = BeautifulSoup(sent_4900, 'lxml')
20 text = soup.get_text()
21 print(text)
```

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

=====

this is yummy, easy and unusual. it makes a quick, delicious pie, crisp or cobbler. home made is better, but a heck of a lot more work. this is great to have on hand for last minute dessert needs where you really want to impress with your creativity in cooking! recommended.

=====

Great flavor, low in calories, high in nutrients, high in protein! Usually protein powders are high priced and high in calories, this one is a great bargain and tastes great, I highly recommend for the lady gym rats, probably not "macho" enough for guys since it is soy based...

=====

For those of you wanting a high-quality, yet affordable green tea, you should definitely give this one a try. Let me first start by saying that everyone is looking for something different for their ideal tea, and I will attempt to briefly highlight what makes this tea attractive to a wide range of tea drinkers (whether you are a beginner or long-time tea enthusiast). I have gone through over 12 boxes of this tea myself, and highly recommend it for the following reasons:
 -Quality: First, this tea offers a smooth quality without any harsh or bitter after tones, which often turns people off from many green teas. I've found my ideal brewing time to be between 3-5 minutes, giving you a light but flavorful cup of tea. However, if you get distracted or forget about your tea and leave it brewing for 20+ minutes like I sometimes do, the quality of this tea is such that you still get a smooth but deeper flavor without the bad after taste. The leaves themselves are whole leaves (not powdered stems, branches, etc commonly found in other brands), and the high-quality nylon bags also include chunks of tropical fruit and other discernible ingredients. This isn't your standard cheap paper bag with a mix of unknown ingredients that have been ground down to a fine powder, leaving you to wonder what it is you are actually drinking.
 -Taste:

This tea offers notes of real pineapple and other hints of tropical fruits, yet isn't sweet or artificially flavored. You have the foundation of a high-quality

lity young hyson green tea for those true "tea flavor" lovers, yet the subtle hints of fruit make this a truly unique tea that I believe most will enjoy. If you want it sweet, you can add sugar, splenda, etc but this really is not necessary as this tea offers an inherent warmth of flavor through it's ingredients.-
 Price: This tea offers an excellent product at an exceptional price (especially when purchased at the prices Amazon offers). Compared to other brands which I believe to be of similar quality (Mighty Leaf, Rishi, Two Leaves, etc.), Revolution offers a superior product at an outstanding price. I have been purchasing this through Amazon for less per box than I would be paying at my local grocery store for Lipton, etc. Overall, this is a wonderful tea that is comparable, and even better than, other teas that are priced much higher. It offers a well-balanced cup of green tea that I believe many will enjoy. In terms of taste, quality, and price, I would argue you won't find a better combination that that offered by Revolution's Tropical Green Tea.

```
In [18]: 1 # https://stackoverflow.com/a/47091490/4084039
2 import re
3
4 def decontracted(phrase):
5     # specific
6     phrase = re.sub(r"won't", "will not", phrase)
7     phrase = re.sub(r"can't", "can not", phrase)
8
9     # general
10    phrase = re.sub(r"n't", " not", phrase)
11    phrase = re.sub(r"\'re", " are", phrase)
12    phrase = re.sub(r"\s", " is", phrase)
13    phrase = re.sub(r"\d", " would", phrase)
14    phrase = re.sub(r"\ll", " will", phrase)
15    phrase = re.sub(r"\t", " not", phrase)
16    phrase = re.sub(r"\ve", " have", phrase)
17    phrase = re.sub(r"\m", " am", phrase)
18    return phrase
```

```
In [19]: 1 sent_1500 = decontracted(sent_1500)
2 print(sent_1500)
3 print("="*50)
```

Great flavor, low in calories, high in nutrients, high in protein! Usually protein powders are high priced and high in calories, this one is a great bargain and tastes great, I highly recommend for the lady gym rats, probably not "macho" enough for guys since it is soy based...

=====

```
In [20]: 1 #remove words with numbers python: https://stackoverflow.com/a/18082370/40840
2 sent_0 = re.sub(r"\S*\d\S*", "", sent_0).strip()
3 print(sent_0)
```

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

```
In [21]: 1 #remove spacial character: https://stackoverflow.com/a/5843547/4084039
2 sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)
3 print(sent_1500)
```

Great flavor low in calories high in nutrients high in protein Usually protein powders are high priced and high in calories this one is a great bargain and tastes great I highly recommend for the lady gym rats probably not macho enough for guys since it is soy based

```
In [22]: 1 # https://gist.github.com/sebleier/554280
2 # we are removing the words from the stop words list: 'no', 'nor', 'not'
3 # <br /><br /> ==> after the above steps, we are getting "br br"
4 # we are including them into stop words list
5 # instead of <br /> if we have <br/> these tags would have revmoved in the 1s
6
7 stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours',
8               "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he',
9               'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'i',
10              'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that',
11              'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have',
12              'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because',
13              'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into',
14              'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on',
15              'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all',
16              'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than',
17              's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "shouldn't",
18              've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn',
19              "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'mustn't',
20              'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't",
21              'won', "won't", 'wouldn', "wouldn't"])
```

```
In [23]: 1 from tqdm import tqdm
2 preprocessed_reviews = []
3 # tqdm is for printing the status bar
4 for sentence in tqdm(final['Text'].values):
5     sentence = re.sub(r"http\S+", "", sentence)
6     sentence = BeautifulSoup(sentence, 'lxml').get_text()
7     sentence = decontracted(sentence)
8     sentence = re.sub("\S*\d\S*", "", sentence).strip()
9     sentence = re.sub('[^A-Za-z]+', ' ', sentence)
10    # https://gist.github.com/sebleier/554280
11    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not in stopwords)
12    preprocessed_reviews.append(sentence.strip())
```

100%|██████████| 46071/46071 [00:27<00:00, 1657.12it/s]

```
In [24]: 1 preprocessed_reviews[1500]
```

```
Out[24]: 'great flavor low calories high nutrients high protein usually protein powders
high priced high calories one great bargain tastes great highly recommend lady
gym rats probably not macho enough guys since soy based'
```

[3.2] Preprocessing Review Summary

```
In [25]: 1 from tqdm import tqdm
2 preprocessed_summary = []
3 # tqdm is for printing the status bar
4 for sentence in tqdm(final['Summary'].values):
5     sentence = re.sub(r"http\S+", "", sentence)
6     sentence = BeautifulSoup(sentence, 'lxml').get_text()
7     sentence = decontracted(sentence)
8     sentence = re.sub("\S*\d\S*", "", sentence).strip()
9     sentence = re.sub('[^A-Za-z]+', ' ', sentence)
10    # https://gist.github.com/sebleier/554280
11    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not
12    preprocessed_summary.append(sentence.strip())
```

76%|██████████ | 34840/46071 [00:15<00:04, 2483.17it/s]C:\Users\sujpanda\Anaconda3\lib\site-packages\bs4__init__.py:219: UserWarning: "b'...'" looks like a filename, not markup. You should probably open this file and pass the filehandle into BeautifulSoup.

'Beautiful Soup.' % markup)

100%|██████████ | 46071/46071 [00:20<00:00, 2221.09it/s]

[5] Assignment 7: SVM

1. Apply SVM 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. Procedure

- You need to work with 2 versions of SVM
 - Linear kernel
 - RBF kernel
- When you are working with linear kernel, use 'SGDClassifier' with hinge loss because it is computationally less expensive.
- When you are working with 'SGDClassifier' with hinge loss and trying to find the AUC score, you would have to use [CalibratedClassifierCV \(https://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html)
- Similarly, like kdtree of knn, when you are working with RBF kernel it's better to reduce the number of dimensions. You can put min_df = 10, max_features = 500 and consider a sample size of 40k points.

3. Hyper parameter tuning (find best alpha in range [10⁻⁴ to 10⁴], and the best penalty among 'l1', 'l2')

- Find the best hyper parameter which will give the maximum [AUC \(https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/\)](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) value

- 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

4. Feature importance


- When you are working on the linear kernel with BOW or TFIDF please print the top 10 best features for each of the positive and negative classes.


5. 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.

6. Representation of results

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

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

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

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

7. Conclusion

- You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this prettytable library [link](http://zetcode.com/python/prettytable/) (<http://zetcode.com/python/prettytable/>)



Note: Data Leakage

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 you train data, and apply the method `transform()` on cv/test data.
4. For more details please go through this [link](https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf). (<https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf>)

Some utility functions

```

In [26]: 1 def check_trade_off(X_train,X_test,y_train,y_test):
2
3     from sklearn.metrics import roc_curve
4     from sklearn.metrics import roc_auc_score
5
6     [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
7
8     C_range1 = ['0.00001', '0.001', '1', '100', '10000',]
9     C_range = [0.00001, 0.001, 1, 100, 10000]
10    dummy_range = [1, 2, 3, 4, 5]
11
12    auc_scores = []
13    auc_train_scores = []
14
15    i = 0
16    for i in C_range:
17        clf = SVC(C=i)
18
19        # fitting the model on crossvalidation train
20        clf.fit(X_train, y_train)
21
22        model = CalibratedClassifierCV(clf, cv='prefit')
23        model.fit(X_train, y_train)
24
25
26        #evaluate AUC score.
27        probs = model.predict_proba(X_test)
28        probs = probs[:, 1]
29        # calculate AUC
30        auc = roc_auc_score(y_test, probs)
31        print('AUC: %.3f' % auc)
32        auc_scores.append(auc)
33
34    print('#####')
35    print('AUC from train data #####')
36    i = 0
37    for i in C_range:
38        clf = SVC(C=i)
39
40        # fitting the model on crossvalidation train
41        clf.fit(X_train, y_train)
42
43        model = CalibratedClassifierCV(clf, cv='prefit')
44        model.fit(X_train, y_train)
45
46        #evaluate AUC score.
47        probs = model.predict_proba(X_train)
48        probs = probs[:, 1]
49        # calculate AUC
50        auc = roc_auc_score(y_train, probs)
51        print('AUC: %.3f' % auc)
52        auc_train_scores.append(auc)
53
54    plt.plot(dummy_range, auc_scores, 'r')
55    plt.plot(dummy_range, auc_train_scores, 'b')
56    plt.xticks(dummy_range, C_range1, rotation='vertical')

```



```
57     for xy in zip(dummy_range, auc_scores):
58         plt.annotate('%f, %f' % xy, xy=xy, textcoords='data')
59     for xy in zip(dummy_range, auc_train_scores):
60         plt.annotate('%f, %f' % xy, xy=xy, textcoords='data')
61
62
63     plt.xlabel('alpha-Values')
64     plt.ylabel('auc_scores')
65     plt.show()
66
```

```

In [27]: 1 def check_trade_off_sgd(X_train,X_test,y_train,y_test):
2
3     from sklearn.metrics import roc_curve
4     from sklearn.metrics import roc_auc_score
5
6     [{'alpha': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
7
8     C_range1 = ['0.00001', '0.001', '1', '100', '10000',]
9     C_range = [0.00001, 0.001, 1, 100, 10000]
10    dummy_range = [1, 2, 3, 4, 5]
11
12    auc_scores = []
13    auc_train_scores = []
14
15    i = 0
16    for i in C_range:
17        clf = SGDClassifier(loss='hinge', alpha=i)
18
19
20        # fitting the model on crossvalidation train
21        clf.fit(X_train, y_train)
22
23        model = CalibratedClassifierCV(clf, cv='prefit')
24        model.fit(X_train, y_train)
25
26
27        #evaluate AUC score.
28        probs = model.predict_proba(X_test)
29        probs = probs[:, 1]
30        # calculate AUC
31        auc = roc_auc_score(y_test, probs)
32        print('AUC: %.3f' % auc)
33        auc_scores.append(auc)
34
35    print('#####')
36    print('AUC from train data #####')
37    i = 0
38    for i in C_range:
39        clf = SGDClassifier(loss='hinge', alpha=i)
40
41
42        # fitting the model on crossvalidation train
43        clf.fit(X_train, y_train)
44
45        model = CalibratedClassifierCV(clf, cv='prefit')
46        model.fit(X_train, y_train)
47
48        #evaluate AUC score.
49        probs = model.predict_proba(X_train)
50        probs = probs[:, 1]
51        # calculate AUC
52        auc = roc_auc_score(y_train, probs)
53        print('AUC: %.3f' % auc)
54        auc_train_scores.append(auc)
55
56    plt.plot(dummy_range, auc_scores, 'r')
57    plt.plot(dummy_range, auc_train_scores, 'b')

```

```

57 plt.xticks(dummy_range, C_range1, rotation='vertical')
58 for xy in zip(dummy_range, auc_scores):
59     plt.annotate('(%f, %f)' % xy, xy=xy, textcoords='data')
60 for xy in zip(dummy_range, auc_train_scores):
61     plt.annotate('(%f, %f)' % xy, xy=xy, textcoords='data')
62
63
64 plt.xlabel('C-Values')
65 plt.ylabel('auc_scores')
66 plt.show()
67

```

In [28]:

```

1 def sgd_results(ialpha,input_penalty,X_train,X_test,y_train,y_test):
2     # roc curve and auc
3     from sklearn.metrics import roc_curve
4     from sklearn.metrics import roc_auc_score
5     from matplotlib import pyplot
6     # ===== Logistic regression=====
7     clf = SGDClassifier(alpha=ialpha,loss='hinge', penalty=input_penalty)
8
9     # fitting the model
10    clf.fit(X_train, y_train)
11
12    # predict the response
13    pred = clf.predict(X_test)
14
15    # evaluate accuracy
16    acc = accuracy_score(y_test, pred) * 100
17    print('\nThe accuracy of the SGD classifier for alpha = %f is %f%%' % (ia
18
19    model = CalibratedClassifierCV(clf, cv='prefit')
20    model.fit(X_train, y_train)
21    probs = model.predict_proba(X_test)
22    probs = probs[:, 1]
23    # calculate AUC
24    auc = roc_auc_score(y_test, probs)
25    print('AUC: %.3f' % auc)
26
27
28    # calculate roc curve
29    fpr, tpr, thresholds = roc_curve(y_test, probs)
30    # plot no skill
31    pyplot.plot([0, 1], [0, 1], linestyle='--')
32    # plot the roc curve for the model
33    pyplot.plot(fpr, tpr, marker='.')
34    # show the plot
35    pyplot.show()
36    from sklearn.metrics import confusion_matrix
37    con_mat = confusion_matrix(y_test, pred, [0, 1])
38    return con_mat,clf

```

```

In [29]: 1 def svc_results(c,X_train,X_test,y_train,y_test):
2         # roc curve and auc
3         from sklearn.metrics import roc_curve
4         from sklearn.metrics import roc_auc_score
5         from matplotlib import pyplot
6         # ===== Logistic regression=====
7         clf = SVC(C=c)
8
9         # fitting the model
10        clf.fit(X_train, y_train)
11
12        # predict the response
13        pred = clf.predict(X_test)
14
15        # evaluate accuracy
16        acc = accuracy_score(y_test, pred) * 100
17        print('\nThe accuracy of the SVC classifier for alpha = %f is %f%%' % (c,
18
19        model = CalibratedClassifierCV(clf, cv='prefit')
20        model.fit(X_train, y_train)
21        probs = model.predict_proba(X_test)
22        probs = probs[:, 1]
23        # calculate AUC
24        auc = roc_auc_score(y_test, probs)
25        print('AUC: %.3f' % auc)
26
27
28        # calculate roc curve
29        fpr, tpr, thresholds = roc_curve(y_test, probs)
30        # plot no skill
31        pyplot.plot([0, 1], [0, 1], linestyle='--')
32        # plot the roc curve for the model
33        pyplot.plot(fpr, tpr, marker='.')
34        # show the plot
35        pyplot.show()
36        from sklearn.metrics import confusion_matrix
37        con_mat = confusion_matrix(y_test, pred, [0, 1])
38        return con_mat,clf

```

```

In [30]: 1 def showHeatMap(con_mat):
2         class_label = ["negative", "positive"]
3         df_cm = pd.DataFrame(con_mat, index = class_label, columns = class_label)
4         sns.heatmap(df_cm, annot = True, fmt = "d")
5         plt.title("Confusion Matrix")
6         plt.xlabel("Predicted Label")
7         plt.ylabel("True Label")
8         plt.show()

```


Applying SVM

[5.1] Linear SVM

[5.1.1] Applying Linear SVM on BOW, SET 1

```
In [246]: 1 from sklearn.cross_validation import train_test_split
2 from sklearn.svm import SVC
3 from sklearn.metrics import accuracy_score
4 from sklearn.cross_validation import cross_val_score
5 from collections import Counter
6 from sklearn.metrics import accuracy_score
7 from sklearn import cross_validation
8 from sklearn.grid_search import GridSearchCV
9 from sklearn.calibration import CalibratedClassifierCV
10 from sklearn.linear_model import SGDClassifier
11 import warnings
12 warnings.filterwarnings("ignore")
```

```
In [247]: 1 X_1, X_test, y_1, y_test = cross_validation.train_test_split(preprocessed_rev
```



In [248]:

```

1 count_vect = CountVectorizer()
2 final_counts = count_vect.fit_transform(X_1)
3 final_test_count = count_vect.transform(X_test)
4
5 # split the train data set into cross validation train and cross validation t
6 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_1, y_1, test_siz
7
8 final_counts_tr_cv = count_vect.transform(X_tr)
9 final_test_count_cv = count_vect.transform(X_cv)
10
11 sgd = SGDClassifier(loss='hinge', penalty='l1', max_iter=1E6, tol=1E-6, shuff
12 tuned_parameters = [{'alpha': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
13
14 sgd_gs = GridSearchCV(sgd, tuned_parameters, scoring='roc_auc', n_jobs=4, cv=
15 sgd_gs.fit(final_counts_tr_cv, y_tr)
16
17 print(sgd_gs.best_estimator_)
18 print(sgd_gs.score(final_test_count_cv, y_cv))
19
20 check_trade_off_sgd(final_counts_tr_cv, final_test_count_cv, y_tr, y_cv)

```

```

SGDClassifier(alpha=0.0001, average=False, class_weight=None, epsilon=0.1,
               eta0=0.0, fit_intercept=True, l1_ratio=0.15,
               learning_rate='optimal', loss='hinge', max_iter=1000000.0,
               n_iter=None, n_jobs=1, penalty='l1', power_t=0.5,
               random_state=123456, shuffle=True, tol=1e-06, verbose=0,
               warm_start=False)

```

0.896747045245148

AUC: 0.907

AUC: 0.925

AUC: 0.511

AUC: 0.581

AUC: 0.581

#####

AUC from train data #####

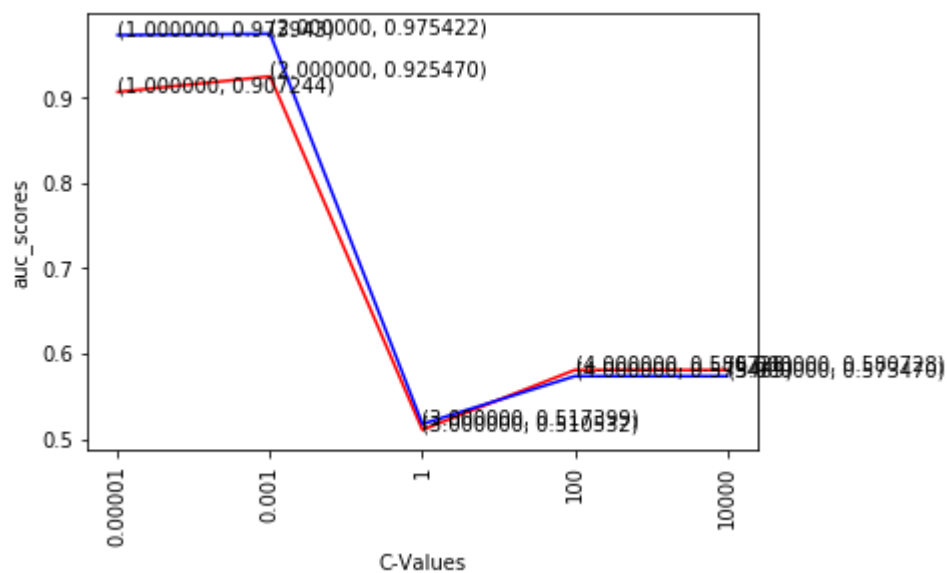
AUC: 0.974

AUC: 0.975

AUC: 0.517

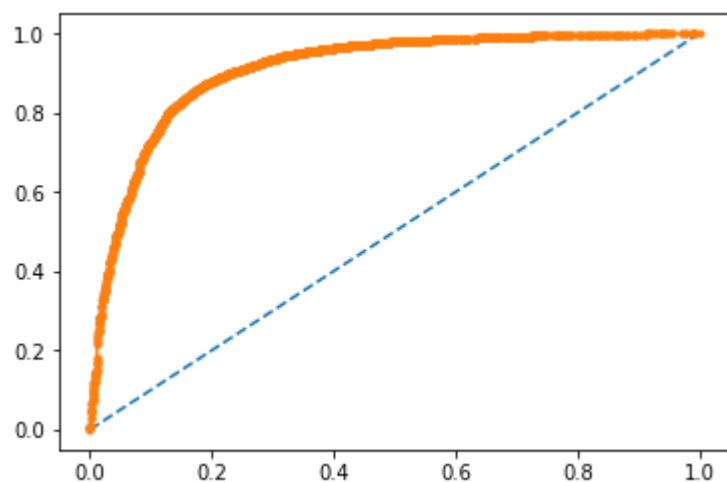
AUC: 0.573

AUC: 0.573



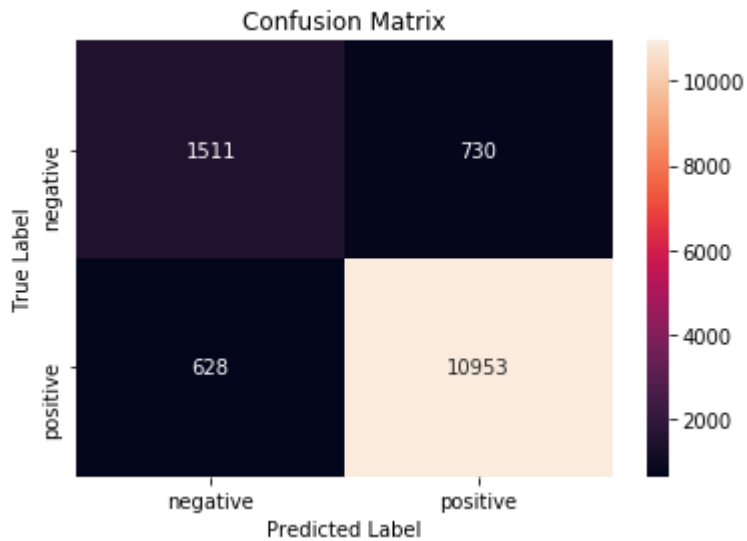
In [249]: 1 con_mat,clf = sgd_results(0.0001,'l1',final_counts,final_test_count,y_1,y_tes

The accuracy of the SGD classifier for alpha = 0.000100 is 90.175083%
AUC: 0.909



Observation: With alpha = 0.0001 the model predicted with accuracy 90% with auc value 0.909.

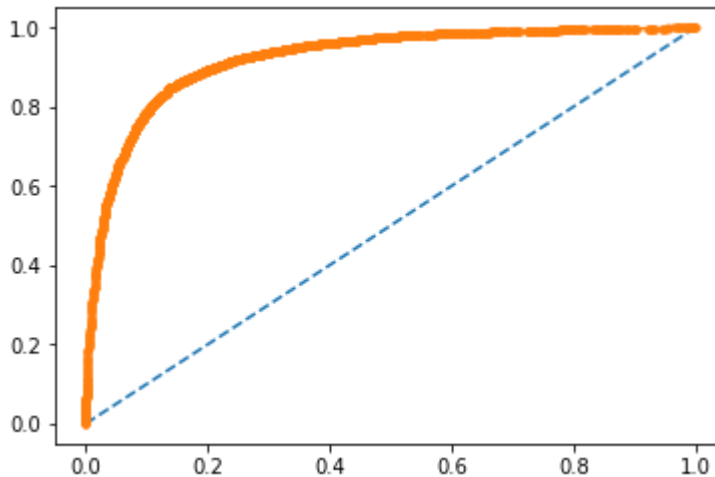
In [250]: 1 showHeatMap(con_mat)



Observation: My model has predicted 628 + 730 points wrongly

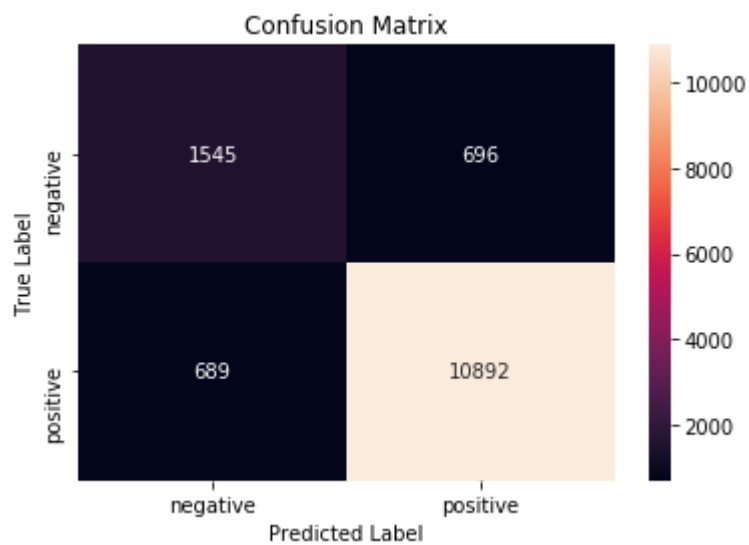
In [251]: 1 con_mat,clf = sgd_results(0.0001,'l2',final_counts,final_test_count,y_1,y_tes

The accuracy of the SGD classifier for alpha = 0.000100 is 89.979742%
AUC: 0.922



Observation: With l2 regularier there is an improvement in AUC score: 0.922 but accuracy reduced to 89. So we can take l2 regularizer for our modeling penalty


```
In [252]: 1 showHeatMap(con_mat)
```



Observation: My model has predicted 689 + 696 points wrongly

[5.1.2] Applying Linear SVM on TFIDF, SET 2

```
In [256]: 1 X_1, X_test, y_1, y_test = cross_validation.train_test_split(preprocessed_rev
```

```
In [257]: 1 tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
2 tf_idf_vect.fit(X_1)
3 final_tf_idf = tf_idf_vect.transform(X_1)
4 final_test_count = tf_idf_vect.transform(X_test)
5
6 # split the train data set into cross validation train and cross validation test
7 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_1, y_1, test_size=0.2)
8
9 final_counts_tr_cv = tf_idf_vect.transform(X_tr)
10 final_test_count_cv = tf_idf_vect.transform(X_cv)
11
12 sgd = SGDClassifier(loss='hinge', penalty='l1')
13 tuned_parameters = [{'alpha': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
14
15 #Using GridSearchCV
16 model = GridSearchCV(sgd, tuned_parameters, scoring = 'roc_auc', cv=5)
17 model.fit(final_counts_tr_cv, y_tr)
18
19 print(model.best_estimator_)
20 print(model.score(final_test_count_cv, y_cv))
21
22 check_trade_off_sgd(final_counts_tr_cv, final_test_count_cv, y_tr, y_cv)
```

```
SGDClassifier(alpha=0.0001, average=False, class_weight=None, epsilon=0.1,
              eta0=0.0, fit_intercept=True, l1_ratio=0.15,
              learning_rate='optimal', loss='hinge', max_iter=None, n_iter=None,
              n_jobs=1, penalty='l1', power_t=0.5, random_state=None,
              shuffle=True, tol=None, verbose=0, warm_start=False)
```

```
0.942661738682314
```

```
AUC: 0.948
```

```
AUC: 0.953
```

```
AUC: 0.608
```

```
AUC: 0.608
```

```
AUC: 0.608
```

```
#####
```

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

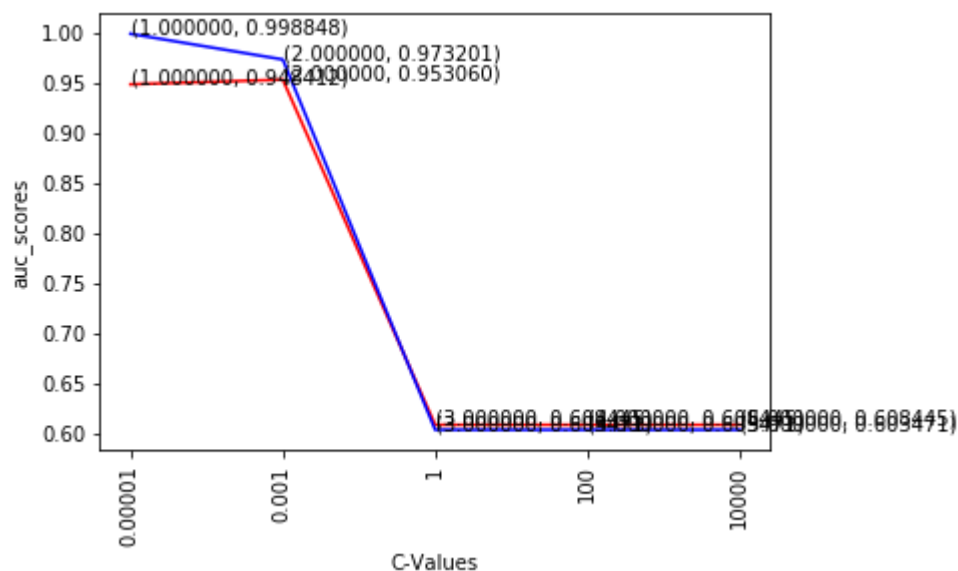
```
AUC: 0.999
```

```
AUC: 0.973
```

```
AUC: 0.603
```

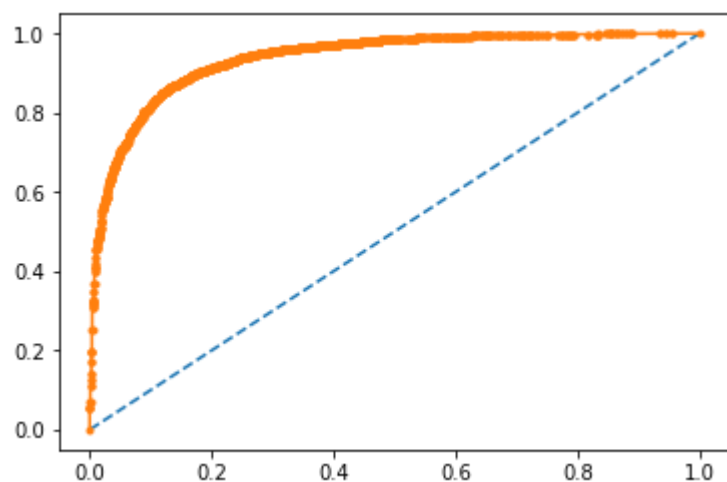
```
AUC: 0.603
```

```
AUC: 0.603
```



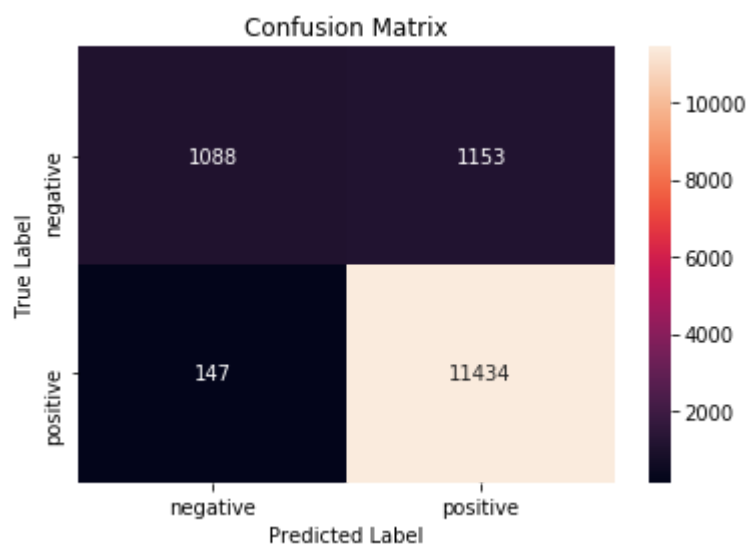
```
In [258]: 1 con_mat,clf = sgd_results(0.0001,'l1',final_tf_idf,final_test_count,y_1,y_tes
```

The accuracy of the SGD classifier for alpha = 0.000100 is 90.594704%
AUC: 0.938



Observation: Model accuracy is 90% with AUC 0.938 when alpha is 0.0001 is very good.

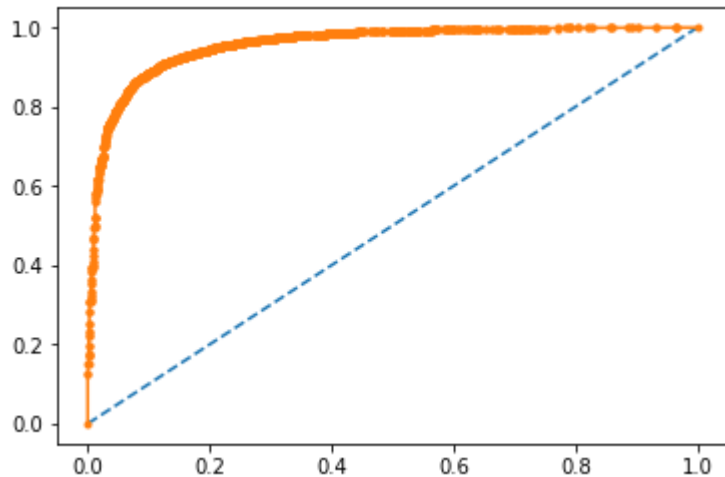
```
In [259]: 1 showHeatMap(con_mat)
```



Observation: My model predicted 147 + 1153 points wrongly

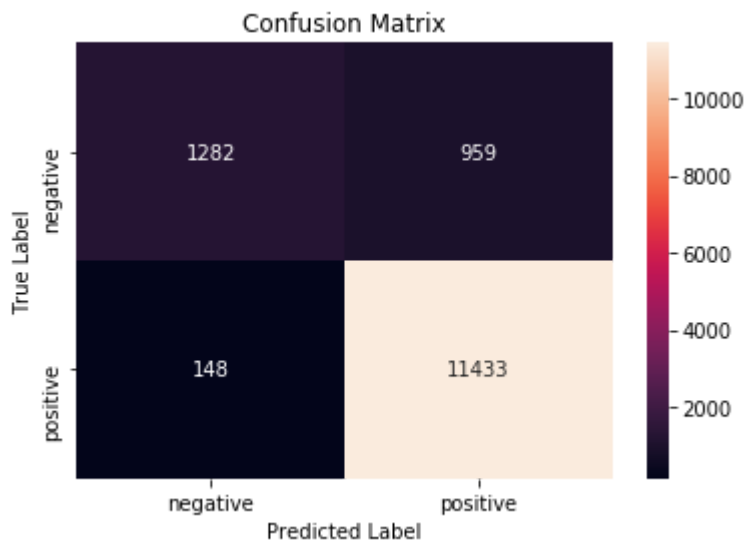
```
In [260]: 1 con_mat,clf = sgdc_results(0.0001,'l2',final_tf_idf,final_test_count,y_1,y_tes
```

The accuracy of the SGD classifier for alpha = 0.000100 is 91.991029%
AUC: 0.957



Observation: There is an improvement in accuracy and AUC score with l2 regularizer

In [261]: 1 showHeatMap(con_mat)



Observation: My model pedicted 148 + 959 points wrongly

In [262]: 1 feature_names = np.array(tf_idf_vect.get_feature_names())
2 sorted_coef_index = clf.coef_[0].argsort()

In [263]: 1 print('Top 10 positive features: \n{}\n'.format(feature_names[sorted_coef_ind
Top 10 positive features:
['great' 'best' 'not disappointed' 'good' 'delicious' 'love' 'loves'
'perfect' 'nice' 'amazing']

In [264]: 1 print('Top 10 negative features: \n{}\n'.format(feature_names[sorted_coef_ind
Top 10 negative features:
['disappointed' 'worst' 'disappointing' 'awful' 'terrible' 'not'
'not worth' 'horrible' 'return' 'not buy']

[5.1.3] Applying Linear SVM on AVG W2V, SET 3

In [265]: 1 X_train, X_test, y_1, y_test = cross_validation.train_test_split(preprocessed

In [266]:

```

1 i=0
2 list_of_sentence=[]
3 for sentence in X_train:
4     list_of_sentence.append(sentence.split())
5
6 w2v_model=Word2Vec(list_of_sentence,min_count=5,size=50, workers=4)
7 w2v_words = list(w2v_model.wv.vocab)
8
9 # average Word2Vec
10 # compute average word2vec for each review.
11 sent_vectors = []; # the avg-w2v for each sentence/review is stored in this L
12 for sent in tqdm(list_of_sentence): # for each review/sentence
13     sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might
14     cnt_words = 0; # num of words with a valid vector in the sentence/review
15     for word in sent: # for each word in a review/sentence
16         if word in w2v_words:
17             vec = w2v_model.wv[word]
18             sent_vec += vec
19             cnt_words += 1
20     if cnt_words != 0:
21         sent_vec /= cnt_words
22     sent_vectors.append(sent_vec)
23 print(sent_vectors[0])
24
25
26
27 i=0
28 list_of_test_sentence=[]
29 for sentence in X_test:
30     list_of_test_sentence.append(sentence.split())
31
32 test_sent_vectors = [];
33
34 for sent in tqdm(list_of_test_sentence): # for each review/sentence
35     sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might
36     cnt_words = 0; # num of words with a valid vector in the sentence/review
37     for word in sent: # for each word in a review/sentence
38         if word in w2v_words:
39             vec = w2v_model.wv[word]
40             sent_vec += vec
41             cnt_words += 1
42     if cnt_words != 0:
43         sent_vec /= cnt_words
44     test_sent_vectors.append(sent_vec)
45 print(test_sent_vectors[0])
46

```

100%|██████████| 32249/32249 [01:37<00:00, 331.40it/s]

```

[-0.70691314  0.05065851  0.38964034  0.41985835 -0.15152602 -0.35196786
-0.25724678  0.76638044 -0.39063415  0.23186132  0.01135624  0.72326949
 0.39849825  0.22387371  0.44969949 -0.22087011 -0.46912268 -0.23675963
-0.70967436  0.25106571 -0.11730615 -0.77854294  0.01344337  0.18702311
-0.22319833 -0.32387671  0.09090055  0.61399192  0.11950142 -0.28587672
-0.34326136  0.1365347   0.43636839 -0.04377889  0.65825292  0.7888547
 0.73269657 -1.08797595  0.14471799  0.07504993  0.56114913  0.75004656

```

0.05591206 1.47118579 0.24307241 0.02181004 -0.1158606 -0.53840384
0.17297038 -0.01418458]

100%|██████████| 13822/13822 [00:40<00:00, 337.49it/s]

[-0.08702452 0.87213708 -0.14136428 -0.53127171 0.27896474 -0.46044726
0.17758902 0.23505041 0.20815672 -0.25931918 0.38914516 0.76943186
0.31039679 -0.02028641 0.36017377 0.06350766 -0.85478444 0.46329409
-0.06438248 0.17123638 0.00314233 -0.85083322 0.40200057 0.90913207
-0.05381919 -0.02819401 0.35263796 0.39815592 0.14521132 -0.13482394
0.28653227 0.84173924 -0.24695739 0.25987838 0.51266812 0.19894371
-0.13075611 0.01975332 0.05176046 0.30230368 -0.04269703 1.40780166
-0.44513858 0.26430924 -0.06118238 0.1416377 -0.17230816 -1.21707915
0.19628136 0.49163068]

In [267]:

```

1  # split the train data set into cross validation train and cross validation test
2  X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_train, y_1, test_size=0.2)
3
4  i=0
5  list_of_cv_sentence=[]
6  for sentence in X_tr:
7      list_of_cv_sentence.append(sentence.split())
8
9  cv_train_sent_vectors = [];
10
11 for sent in tqdm(list_of_cv_sentence): # for each review/sentence
12     sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to initialize them
13     cnt_words = 0; # num of words with a valid vector in the sentence/review
14     for word in sent: # for each word in a review/sentence
15         if word in w2v_words:
16             vec = w2v_model.wv[word]
17             sent_vec += vec
18             cnt_words += 1
19     if cnt_words != 0:
20         sent_vec /= cnt_words
21     cv_train_sent_vectors.append(sent_vec)
22 print(cv_train_sent_vectors[0])
23
24 i=0
25 list_of_cv_test_sentence=[]
26 for sentence in X_cv:
27     list_of_cv_test_sentence.append(sentence.split())
28
29 cv_test_sent_vectors = [];
30
31 for sent in tqdm(list_of_cv_test_sentence): # for each review/sentence
32     sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to initialize them
33     cnt_words = 0; # num of words with a valid vector in the sentence/review
34     for word in sent: # for each word in a review/sentence
35         if word in w2v_words:
36             vec = w2v_model.wv[word]
37             sent_vec += vec
38             cnt_words += 1
39     if cnt_words != 0:
40         sent_vec /= cnt_words
41     cv_test_sent_vectors.append(sent_vec)
42 print(cv_test_sent_vectors[0])
43
44 sgd = SGDClassifier(loss='hinge', penalty='l1')
45 tuned_parameters = [{'alpha': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
46
47 #Using GridSearchCV
48 model = GridSearchCV(sgd, tuned_parameters, scoring = 'roc_auc', cv=5)
49 model.fit(cv_train_sent_vectors, y_tr)
50
51 print(model.best_estimator_)
52 print(model.score(cv_test_sent_vectors, y_cv))
53
54 check_trade_off_sgd(cv_train_sent_vectors, cv_test_sent_vectors, y_tr, y_cv)

```

100% |██████████| 22574/22574 [01:04<00:00, 350.10it/s]


```
[ 0.08510495  0.26092592 -0.79238879 -1.3508834  0.1863163  0.11327625
-0.53192805 -0.69178598 -0.49748244  0.04794091  0.24380035  0.4417329
 1.32874974 -0.32551921  0.41531095  1.03252378  0.13753669 -0.7313178
 0.18786763  0.21914234 -0.47548723 -0.63648485  0.80846777  1.36391416
 0.36101585 -0.2344835  0.30594579  0.58087065 -1.37438641 -0.61767826
 0.05150555  0.92075367 -1.19195867  0.28509153  0.54597462 -0.3566633
-0.58560209 -0.17822589 -0.1889485  0.56225488  0.4105128 -0.62780031
-0.40520856  1.39424051  0.89706699  0.42032172 -0.78082825 -0.82834497
-0.70316356 -0.47892136]
```

100%|██████████| 9675/9675 [00:27<00:00, 353.71it/s]

```
[ 0.23674532  0.48092168 -0.48189216 -1.19358868 -0.30842771 -0.07806226
-0.10428937 -0.18510295 -0.00796553 -0.43365363  0.60801747  0.57487001
 0.12955645  0.24113405  0.0834538 -0.33371197 -0.41175283 -0.21911259
 0.06936753  0.14725839 -0.57582733 -0.4676146  0.47448143  0.67059709
 0.15942912 -0.18546413 -0.08690518 -0.05383114 -0.57850386  0.08525441
 0.21540359  0.6891708 -0.79119456  0.54539031  0.28150502 -0.2014107
 0.1465843  0.39709588  0.58634039  0.06259895 -0.2334361 -0.02668138
-0.07403829  0.26700662  0.00827627 -0.24273524  0.08547828 -1.23042495
-0.26745986  0.25093822]
```

```
SGDClassifier(alpha=0.0001, average=False, class_weight=None, epsilon=0.1,
              eta0=0.0, fit_intercept=True, l1_ratio=0.15,
              learning_rate='optimal', loss='hinge', max_iter=None, n_iter=None,
              n_jobs=1, penalty='l1', power_t=0.5, random_state=None,
              shuffle=True, tol=None, verbose=0, warm_start=False)
```

0.8877449843715736

AUC: 0.867

AUC: 0.894

AUC: 0.893

AUC: 0.638

AUC: 0.638

#####

AUC from train data #####

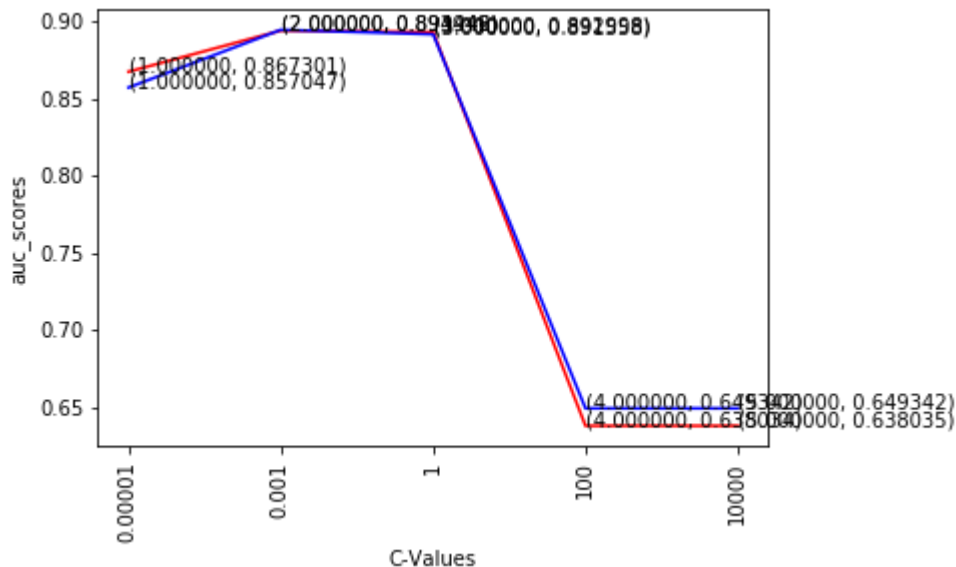
AUC: 0.857

AUC: 0.894

AUC: 0.892

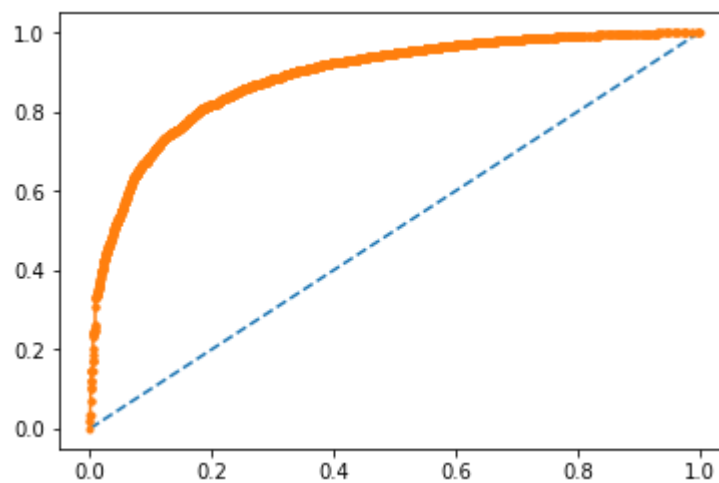
AUC: 0.649

AUC: 0.649



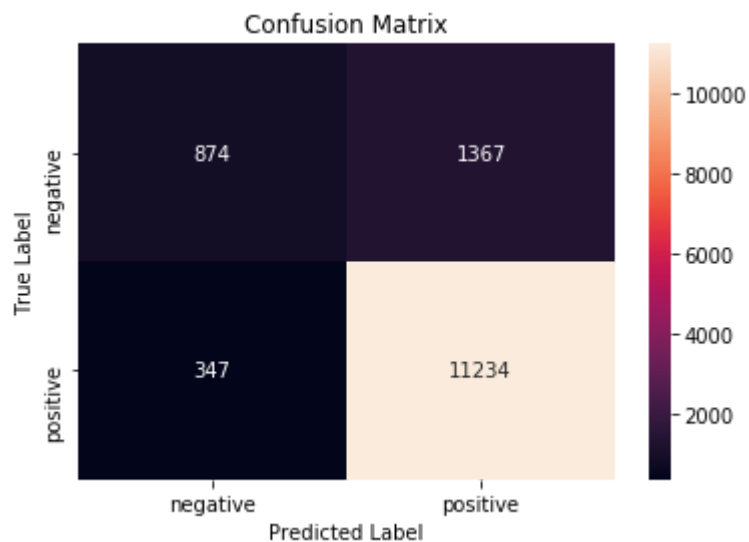
```
In [268]: 1 con_mat,clf = sgd_results(0.0001,'l1',sent_vectors,test_sent_vectors,y_1,y_te
```

The accuracy of the SGD classifier for alpha = 0.000100 is 87.599479%
AUC: 0.888



Observation: My model has predicted with 87 accuracy with AUC: 0.888

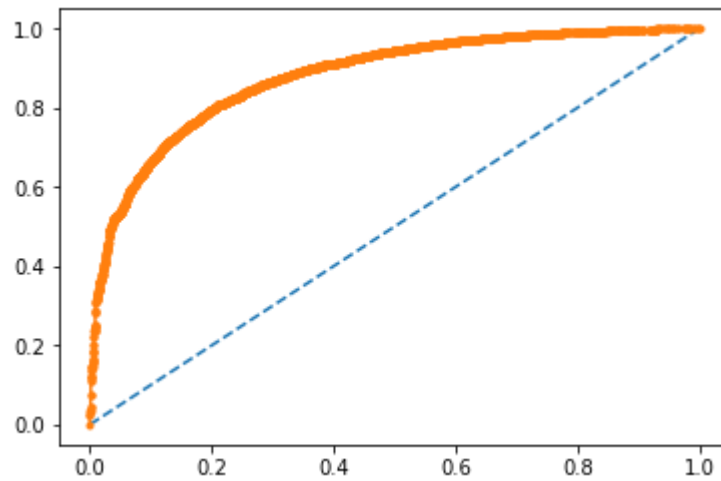
```
In [269]: 1 showHeatMap(con_mat)
```



Observation: My model predicted 347 + 1367 points wrongly

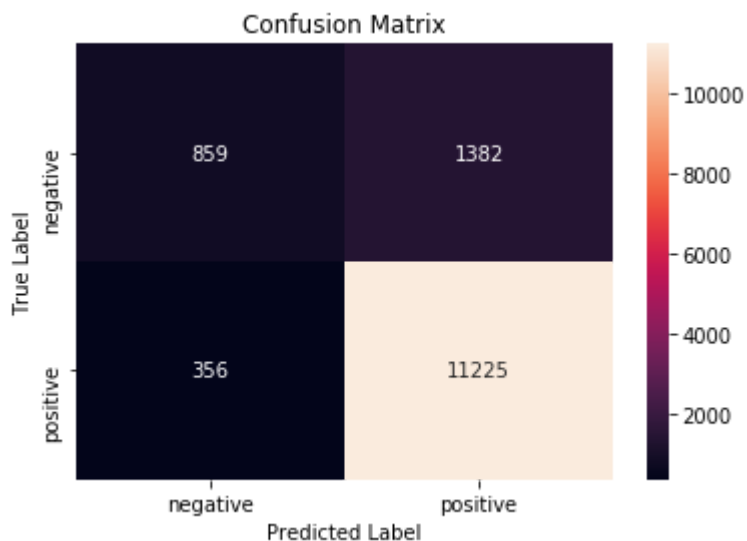
```
In [270]: 1 con_mat,clf = sgd_results(0.0001,'l2',sent_vectors,test_sent_vectors,y_1,y_te
```

The accuracy of the SGD classifier for alpha = 0.000100 is 87.425843%
AUC: 0.880



Observation: Not much difference between l1 and l2 regularizer output

```
In [271]: 1 showHeatMap(con_mat)
```



Observation: My model predicted 356 + 1382 points wrongly

[5.1.4] Applying Linear SVM on TFIDF W2V, SET 4

```
In [272]: 1 X_train, X_test, y_1, y_test = cross_validation.train_test_split(preprocessed
```

```
In [273]: 1 model = TfidfVectorizer()
2 X_train_transformed = model.fit_transform(X_train)
3
4 dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
```

```
In [274]: 1 # Train your own Word2Vec model using your own text corpus
2 i=0
3 list_of_sentence=[]
4 for sentence in X_train:
5     list_of_sentence.append(sentence.split())
```

```
In [275]: 1 w2v_model=Word2Vec(list_of_sentence,min_count=5,size=50, workers=4)
2 w2v_words = list(w2v_model.wv.vocab)
```

```
In [276]: 1 # TF-IDF weighted Word2Vec
2 tfidf_feat = model.get_feature_names() # tfidf words/col-names
3 # final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val
4
5 tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is stored i
6 row=0;
7 for sent in tqdm(list_of_sentence): # for each review/sentence
8     sent_vec = np.zeros(50) # as word vectors are of zero length
9     weight_sum =0; # num of words with a valid vector in the sentence/review
10    for word in sent: # for each word in a review/sentence
11        if word in w2v_words and word in tfidf_feat:
12            vec = w2v_model.wv[word]
13            # tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
14            # to reduce the computation we are
15            # dictionary[word] = idf value of word in whole courpus
16            # sent.count(word) = tf valeus of word in this review
17            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
18            sent_vec += (vec * tf_idf)
19            weight_sum += tf_idf
20    if weight_sum != 0:
21        sent_vec /= weight_sum
22    tfidf_sent_vectors.append(sent_vec)
23    row += 1
```

100%|██████████| 32249/32249 [18:54<00:00, 28.41it/s]

```
In [277]: 1 i=0
2 list_of_test_sentence=[]
3 for sentence in X_test:
4     list_of_test_sentence.append(sentence.split())
```

In [278]:

```

1  # TF-IDF weighted Word2Vec
2  tfidf_feat = model.get_feature_names() # tfidf words/col-names
3  # final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val
4
5  tfidf_test_sent_vectors = []; # the tfidf-w2v for each sentence/review is sto
6  row=0;
7  for sent in tqdm(list_of_test_sentence): # for each review/sentence
8      sent_vec = np.zeros(50) # as word vectors are of zero length
9      weight_sum =0; # num of words with a valid vector in the sentence/review
10     for word in sent: # for each word in a review/sentence
11         if word in w2v_words and word in tfidf_feat:
12             vec = w2v_model.wv[word]
13             # tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
14             # to reduce the computation we are
15             # dictionary[word] = idf value of word in whole corpus
16             # sent.count(word) = tf value of word in this review
17             tf_idf = dictionary[word]*(sent.count(word)/len(sent))
18             sent_vec += (vec * tf_idf)
19             weight_sum += tf_idf
20     if weight_sum != 0:
21         sent_vec /= weight_sum
22     tfidf_test_sent_vectors.append(sent_vec)
23     row += 1

```

100%|██████████| 13822/13822 [08:23<00:00, 27.45it/s]

```

In [279]: 1 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_train, y_1, test
2
3 i=0
4 list_of_cv_sentence=[]
5 for sentence in X_tr:
6     list_of_cv_sentence.append(sentence.split())
7
8
9 tfidf_feat = model.get_feature_names() # tfidf words/col-names
10 # final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val
11
12 tfidf_cv_sent_vectors = []; # the tfidf-w2v for each sentence/review is store
13 row=0;
14 for sent in tqdm(list_of_cv_sentence): # for each review/sentence
15     sent_vec = np.zeros(50) # as word vectors are of zero length
16     weight_sum =0; # num of words with a valid vector in the sentence/review
17     for word in sent: # for each word in a review/sentence
18         if word in w2v_words and word in tfidf_feat:
19             vec = w2v_model.wv[word]
20             # tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
21             # to reduce the computation we are
22             # dictionary[word] = idf value of word in whole courpus
23             # sent.count(word) = tf valeus of word in this review
24             tf_idf = dictionary[word]*(sent.count(word)/len(sent))
25             sent_vec += (vec * tf_idf)
26             weight_sum += tf_idf
27         if weight_sum != 0:
28             sent_vec /= weight_sum
29         tfidf_cv_sent_vectors.append(sent_vec)
30         row += 1
31
32 i=0
33 list_of_cv_test_sentence=[]
34 for sentence in X_cv:
35     list_of_cv_test_sentence.append(sentence.split())
36
37
38 tfidf_cv_test_sent_vectors = []; # the tfidf-w2v for each sentence/review is
39 row=0;
40 for sent in tqdm(list_of_cv_test_sentence): # for each review/sentence
41     sent_vec = np.zeros(50) # as word vectors are of zero length
42     weight_sum =0; # num of words with a valid vector in the sentence/review
43     for word in sent: # for each word in a review/sentence
44         if word in w2v_words and word in tfidf_feat:
45             vec = w2v_model.wv[word]
46             # tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
47             # to reduce the computation we are
48             # dictionary[word] = idf value of word in whole courpus
49             # sent.count(word) = tf valeus of word in this review
50             tf_idf = dictionary[word]*(sent.count(word)/len(sent))
51             sent_vec += (vec * tf_idf)
52             weight_sum += tf_idf
53         if weight_sum != 0:
54             sent_vec /= weight_sum
55         tfidf_cv_test_sent_vectors.append(sent_vec)
56         row += 1

```

```

57
58
59 sgd = SGDClassifier(loss='hinge', penalty='l1')
60 tuned_parameters = [{'alpha': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
61
62 #Using GridSearchCV
63 model = GridSearchCV(sgd, tuned_parameters, scoring = 'roc_auc', cv=5)
64 model.fit(tfidf_cv_sent_vectors, y_tr)
65
66 print(model.best_estimator_)
67 print(model.score(tfidf_cv_test_sent_vectors, y_cv))
68
69 check_trade_off_sgd(tfidf_cv_sent_vectors,tfidf_cv_test_sent_vectors,y_tr,y_c
70

```

```
100%|██████████| 22574/22574 [13:19<00:00, 28.25it/s]
```

```
100%|██████████| 9675/9675 [06:20<00:00, 25.42it/s]
```

```

SGDClassifier(alpha=0.0001, average=False, class_weight=None, epsilon=0.1,
              eta0=0.0, fit_intercept=True, l1_ratio=0.15,
              learning_rate='optimal', loss='hinge', max_iter=None, n_iter=None,
              n_jobs=1, penalty='l1', power_t=0.5, random_state=None,
              shuffle=True, tol=None, verbose=0, warm_start=False)

```

```
0.8630297111241529
```

```
AUC: 0.828
```

```
AUC: 0.868
```

```
AUC: 0.833
```

```
AUC: 0.634
```

```
AUC: 0.634
```

```
#####
```

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

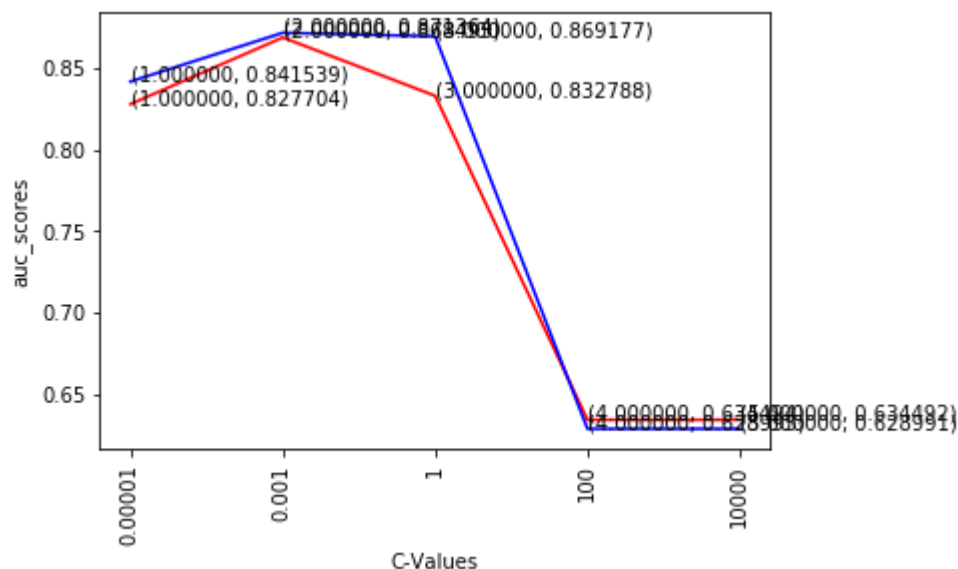
```
AUC: 0.842
```

```
AUC: 0.871
```

```
AUC: 0.869
```

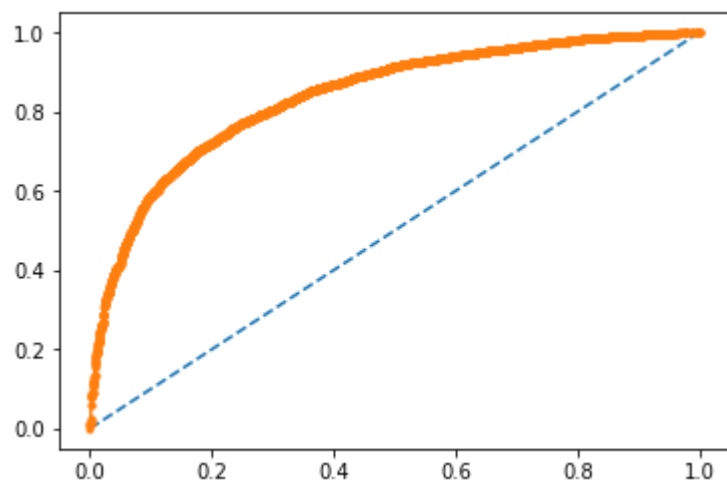
```
AUC: 0.629
```

```
AUC: 0.629
```



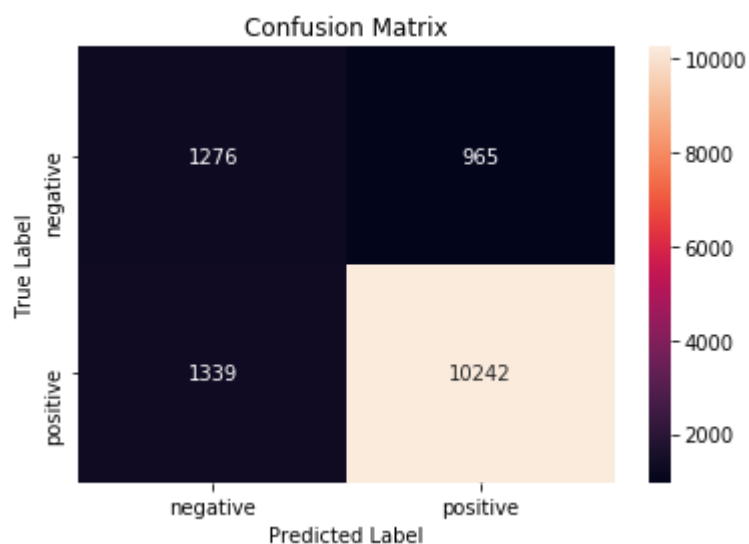
```
In [280]: 1 con_mat,clf = sgdf_results(0.0001,'11',tfidf_sent_vectors,tfidf_test_sent_vect
```

The accuracy of the SGD classifier for alpha = 0.000100 is 83.330922%
AUC: 0.840



Observation: Model predicted with accuracy 83 % and 0.840 auc score.

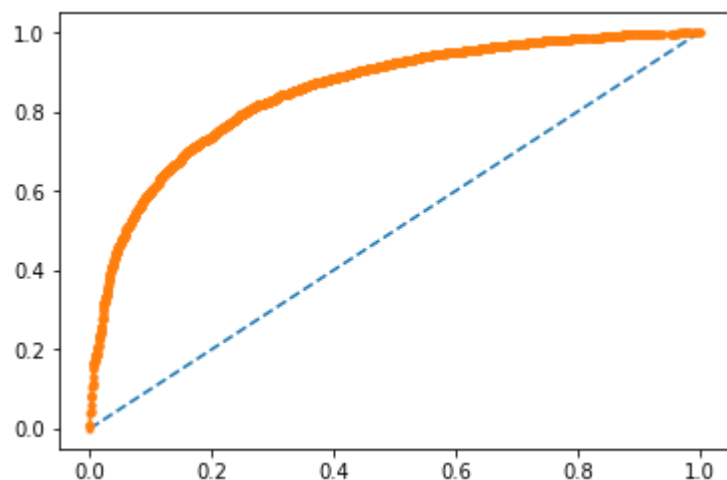
```
In [281]: 1 showHeatMap(con_mat)
```



Observation: My model predicted 1339 + 965 points wrongly

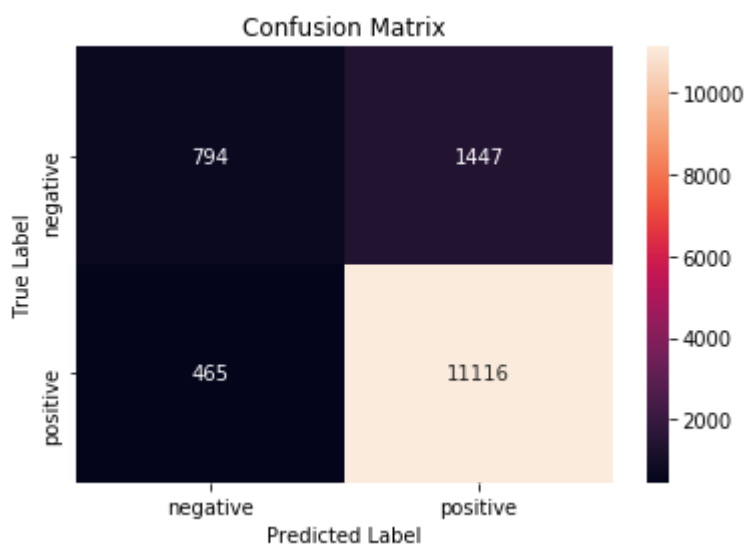

```
In [282]: 1 con_mat,clf = sgd_results(0.0001,'l2',tfidf_sent_vectors,tfidf_test_sent_vect
```

The accuracy of the SGD classifier for alpha = 0.000100 is 86.166980%
AUC: 0.852



Observation: There is not much difference between l1 and l2 regularizer

```
In [283]: 1 showHeatMap(con_mat)
```



Observation: My model predicted 456 + 1447 points wrongly

[5.2] RBF SVM

[5.2.1] Applying RBF SVM on BOW, SET 1

In [284]:

```
1 X_1, X_test, y_1, y_test = cross_validation.train_test_split(preprocessed_rev  
2
```



```

In [285]: 1 count_vect = CountVectorizer(min_df=10, max_features=500)
2 final_counts = count_vect.fit_transform(X_1)
3 final_test_count = count_vect.transform(X_test)
4
5 # split the train data set into cross validation train and cross validation t
6 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_1, y_1, test_siz
7
8 final_counts_tr_cv = count_vect.transform(X_tr)
9 final_test_count_cv = count_vect.transform(X_cv)
10
11
12 tuned_parameters = [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
13
14 #Using GridSearchCV
15 model = GridSearchCV(SVC(), tuned_parameters, scoring = 'roc_auc', cv=5)
16 model.fit(final_counts_tr_cv, y_tr)
17
18 print(model.best_estimator_)
19 print(model.score(final_test_count_cv, y_cv))
20 check_trade_off(final_counts_tr_cv, final_test_count_cv, y_tr, y_cv)

```

```

SVC(C=100, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto', kernel='rbf',
    max_iter=-1, probability=False, random_state=None, shrinking=True,
    tol=0.001, verbose=False)

```

```
0.9029186572977717
```

```
AUC: 0.676
```

```
AUC: 0.885
```

```
AUC: 0.893
```

```
AUC: 0.903
```

```
AUC: 0.847
```

```
#####
```

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

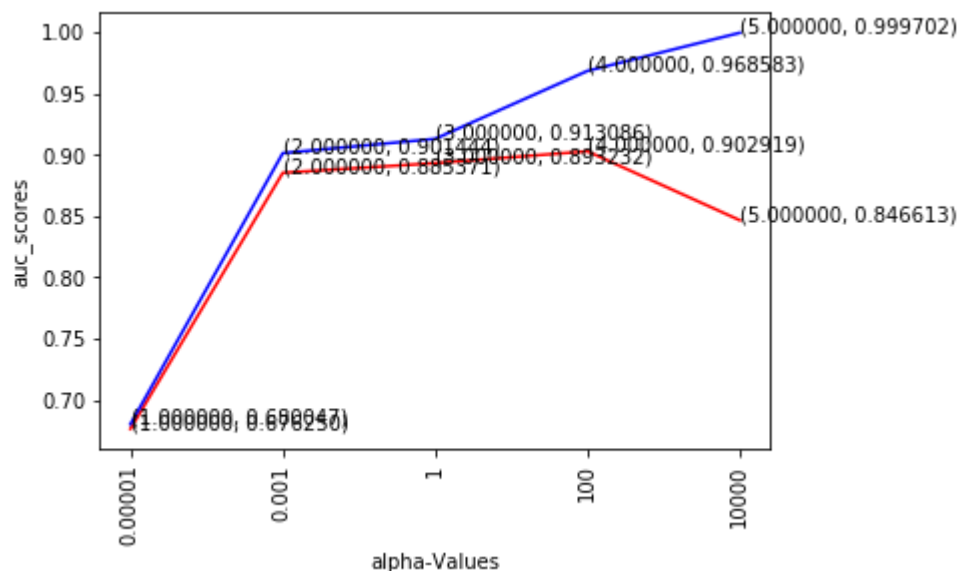
```
AUC: 0.680
```

```
AUC: 0.901
```

```
AUC: 0.913
```

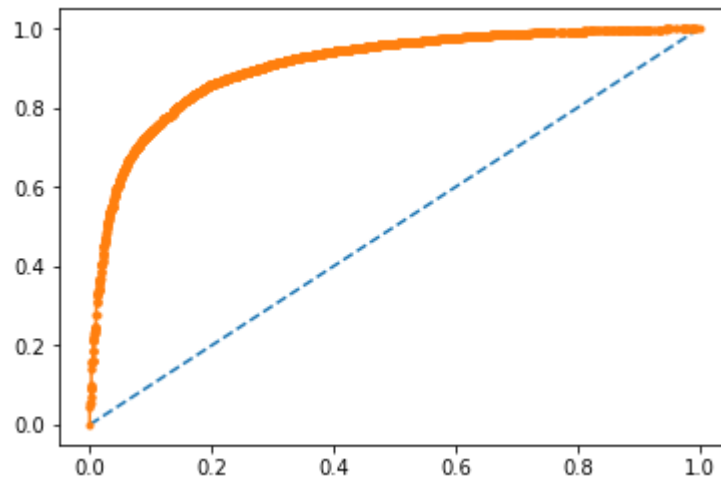
```
AUC: 0.969
```

```
AUC: 1.000
```



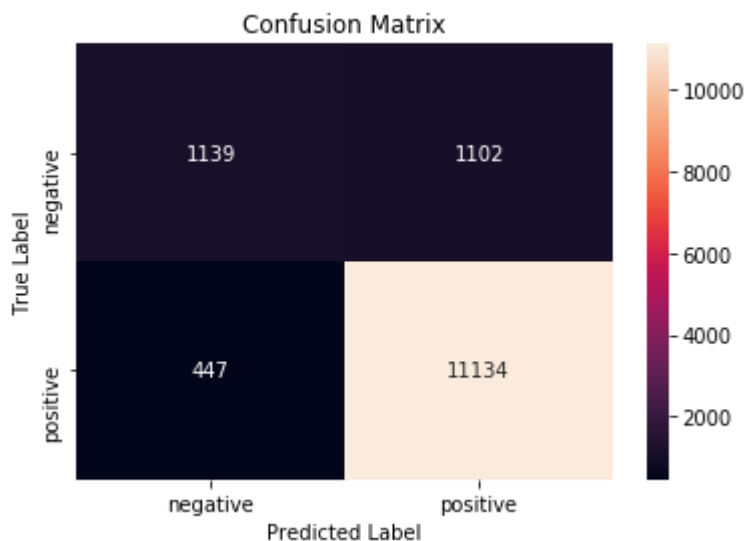
```
In [286]: 1 con_mat,clf = svc_results(100,final_counts,final_test_count,y_1,y_test)
```

The accuracy of the SGD classifier for alpha = 100.000000 is 88.793228%
AUC: 0.906



observation: The model predicted 88% accuracy with AUC: 906 when c = 100

```
In [287]: 1 showHeatMap(con_mat)
```



Observation: This model predicted 447 + 1102 points wrongly

[5.2.2] Applying RBF SVM on TFIDF, SET 2

```
In [288]: 1 X_1, X_test, y_1, y_test = cross_validation.train_test_split(preprocessed_rev
```

```
In [289]: 1 tf_idf_vect = TfidfVectorizer(min_df=10, max_features=500)
2 tf_idf_vect.fit(X_1)
3 final_tf_idf = tf_idf_vect.transform(X_1)
4 final_test_count = tf_idf_vect.transform(X_test)
5
6 # split the train data set into cross validation train and cross validation test
7 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_1, y_1, test_size=0.2)
8
9 final_counts_tr_cv = tf_idf_vect.transform(X_tr)
10 final_test_count_cv = tf_idf_vect.transform(X_cv)
11
12 tuned_parameters = [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
13
14 #Using GridSearchCV
15 model = GridSearchCV(SVC(), tuned_parameters, scoring = 'roc_auc', cv=5)
16 model.fit(final_counts_tr_cv, y_tr)
17
18 print(model.best_estimator_)
19 print(model.score(final_test_count_cv, y_cv))
20
21 check_trade_off(final_counts_tr_cv, final_test_count_cv, y_tr, y_cv)
```

```
SVC(C=10000, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto', kernel='rbf',
    max_iter=-1, probability=False, random_state=None, shrinking=True,
    tol=0.001, verbose=False)
```

```
0.9061422815171389
```

```
AUC: 0.843
```

```
AUC: 0.843
```

```
AUC: 0.900
```

```
AUC: 0.904
```

```
AUC: 0.906
```

```
#####
```

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

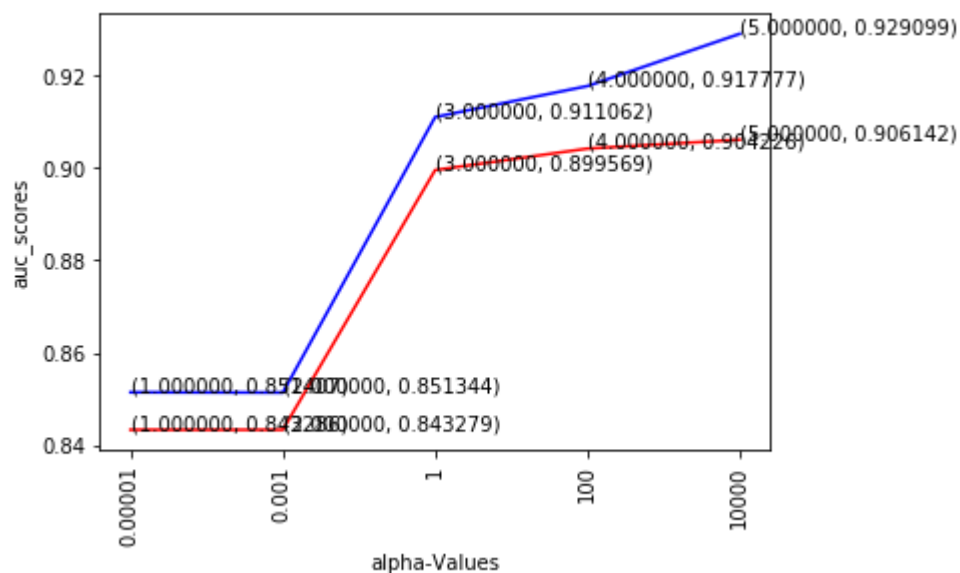
```
AUC: 0.851
```

```
AUC: 0.851
```

```
AUC: 0.911
```

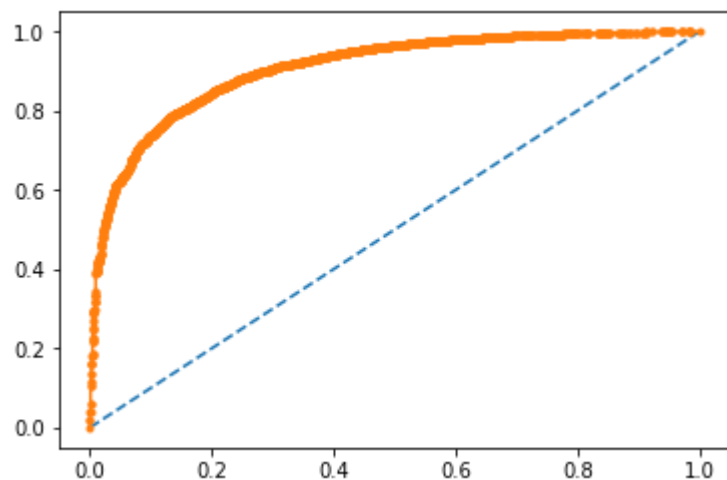
```
AUC: 0.918
```

```
AUC: 0.929
```



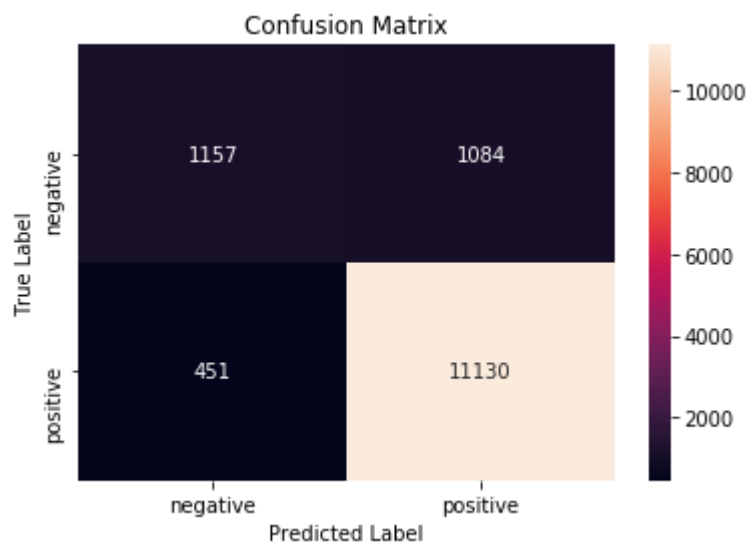
```
In [291]: 1 con_mat,clf = svc_results(10000,final_tf_idf,final_test_count,y_1,y_test)
```

The accuracy of the SVC classifier for alpha = 10000.000000 is 88.894516%
AUC: 0.907



Observation: Model predicted 88% accuracy with AUC score: 0.907

```
In [293]: 1 showHeatMap(con_mat)
```



Observation: My model predicted 451 + 1084 points wrongly

[5.2.3] Applying RBF SVM on AVG W2V, SET 3

```
In [294]: 1 X_1, X_test, y_1, y_test = cross_validation.train_test_split(preprocessed_rev
```

In [295]:

```

1 i=0
2 list_of_sentence=[]
3 for sentence in X_train:
4     list_of_sentence.append(sentence.split())
5
6 w2v_model=Word2Vec(list_of_sentence,min_count=5,size=50, workers=4)
7 w2v_words = list(w2v_model.wv.vocab)
8
9 # average Word2Vec
10 # compute average word2vec for each review.
11 sent_vectors = []; # the avg-w2v for each sentence/review is stored in this list
12 for sent in tqdm(list_of_sentence): # for each review/sentence
13     sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to initialize them
14     cnt_words = 0; # num of words with a valid vector in the sentence/review
15     for word in sent: # for each word in a review/sentence
16         if word in w2v_words:
17             vec = w2v_model.wv[word]
18             sent_vec += vec
19             cnt_words += 1
20     if cnt_words != 0:
21         sent_vec /= cnt_words
22     sent_vectors.append(sent_vec)
23 print(sent_vectors[0])
24
25
26
27 i=0
28 list_of_test_sentence=[]
29 for sentence in X_test:
30     list_of_test_sentence.append(sentence.split())
31
32 test_sent_vectors = [];
33
34 for sent in tqdm(list_of_test_sentence): # for each review/sentence
35     sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to initialize them
36     cnt_words = 0; # num of words with a valid vector in the sentence/review
37     for word in sent: # for each word in a review/sentence
38         if word in w2v_words:
39             vec = w2v_model.wv[word]
40             sent_vec += vec
41             cnt_words += 1
42     if cnt_words != 0:
43         sent_vec /= cnt_words
44     test_sent_vectors.append(sent_vec)
45 print(test_sent_vectors[0])
46

```

100%|██████████| 32249/32249 [01:32<00:00, 348.24it/s]

```

[-0.3229212  0.0909635  0.14200178  0.43835705  0.18845611 -0.43789002
 0.05517153  0.70787672 -0.71328895  0.77850036  0.26874793  0.95412588
 0.52620311  0.11158776  0.10857141  0.09850775 -0.52109222 -0.17406479
-0.85310262  0.08417704 -0.34046573 -0.68380223  0.01221205  0.05715774
-0.26670442 -0.43873761 -0.44164987  0.56933045 -0.13125047 -0.35325936
-0.49677751  0.12477424  0.3830783  0.04410631  0.38815009  0.62609032
 0.45125122 -0.9022373  0.13413235  0.08155093  0.4543668  0.93590224]

```


0.11707443 1.46638633 -0.04939898 0.35705601 0.15852814 -0.51596832
0.25066362 -0.20028094]

100%|██████████| 13822/13822 [00:41<00:00, 332.90it/s]

[-0.10094378 1.11984005 -0.04321185 -0.4043802 0.37045777 -0.63383948
0.09150795 0.10056138 0.37985576 -0.31994779 0.46529526 0.83320616
0.40134423 -0.02841795 0.2435701 0.06110854 -0.9379301 0.53725673
-0.12202664 0.35432463 0.04880635 -0.70764226 0.28782641 1.04118215
0.07696167 0.04706393 -0.05136644 0.38052209 0.03635978 0.04588852
0.077943 0.89141433 -0.13332699 0.31047307 0.57601362 0.16244922
-0.20730958 -0.02952294 -0.17687802 -0.08034172 0.16575303 1.23290972
-0.23976431 0.58490628 -0.08059081 0.6494041 0.17343514 -0.8229991
-0.26132654 0.60701266]

In [296]:

```

1  # split the train data set into cross validation train and cross validation test
2  X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_train, y_1, test_size=0.2)
3
4  i=0
5  list_of_cv_sentence=[]
6  for sentence in X_tr:
7      list_of_cv_sentence.append(sentence.split())
8
9  cv_train_sent_vectors = [];
10
11 for sent in tqdm(list_of_cv_sentence): # for each review/sentence
12     sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to initialize them
13     cnt_words = 0; # num of words with a valid vector in the sentence/review
14     for word in sent: # for each word in a review/sentence
15         if word in w2v_words:
16             vec = w2v_model.wv[word]
17             sent_vec += vec
18             cnt_words += 1
19     if cnt_words != 0:
20         sent_vec /= cnt_words
21     cv_train_sent_vectors.append(sent_vec)
22 print(cv_train_sent_vectors[0])
23
24 i=0
25 list_of_cv_test_sentence=[]
26 for sentence in X_cv:
27     list_of_cv_test_sentence.append(sentence.split())
28
29 cv_test_sent_vectors = [];
30
31 for sent in tqdm(list_of_cv_test_sentence): # for each review/sentence
32     sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to initialize them
33     cnt_words = 0; # num of words with a valid vector in the sentence/review
34     for word in sent: # for each word in a review/sentence
35         if word in w2v_words:
36             vec = w2v_model.wv[word]
37             sent_vec += vec
38             cnt_words += 1
39     if cnt_words != 0:
40         sent_vec /= cnt_words
41     cv_test_sent_vectors.append(sent_vec)
42 print(cv_test_sent_vectors[0])
43
44 tuned_parameters = [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
45
46 #Using GridSearchCV
47 model = GridSearchCV(SVC(), tuned_parameters, scoring = 'roc_auc', cv=5)
48 model.fit(cv_train_sent_vectors, y_tr)
49
50 print(model.best_estimator_)
51 print(model.score(cv_test_sent_vectors, y_cv))
52
53 check_trade_off(cv_train_sent_vectors, cv_test_sent_vectors, y_tr, y_cv)

```

100%|██████████| 22574/22574 [01:05<00:00, 346.54it/s]

```

[-0.2905551  0.62139706 -0.03318877 -0.49136731  0.09006259 -0.01776666
-0.17018281  0.03852796  0.35760078  0.04350535  0.48829077  0.52911945
 0.16666693  0.03080599  0.50328009 -0.29323572 -0.80406415 -0.32414489
-0.02110795 -0.17709955 -0.57182706 -0.35559807  0.11828896  0.55597643
 0.04393607 -0.31599329 -0.05026706  0.42649929 -0.15260916  0.11197962
 0.26053023  0.20524775 -0.2526622  0.29333542  0.07927508 -0.18340579
 0.00891723 -0.29392766  0.20088845  0.13866965 -0.08526035  0.43451062
-0.037705   0.57984516 -0.03342701  0.38694469 -0.0286005  -0.62059804
-0.39589116 -0.13186588]

```

100%|██████████| 9675/9675 [00:28<00:00, 343.38it/s]

```

[-0.73731186  0.45247821 -0.20114915 -0.77948135  0.69915743 -0.60044611
-0.48607216 -0.44931495  0.15578087 -0.24672109  0.19700561  1.08026159
-0.1196873  -0.95607778 -0.01002884  0.3647039  0.19647758 -0.20182659
-0.53337809  0.34062649 -1.00834926 -0.34266816  0.41765108  0.17859725
-0.23061378 -0.64971501  0.1354455  1.22837861 -0.83665973 -0.76216334
-1.07714029 -0.02715902 -0.48696535  0.53525911  0.90515195  0.12492945
-0.33667326  0.88673962  0.09865461  0.91019244 -0.04071382 -0.70085437
-0.0475254  -0.06126182 -0.20639705  0.24520903 -0.72075997 -0.68022136
-0.33296668 -0.33773697]

```

```

SVC(C=100, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto', kernel='rbf',
    max_iter=-1, probability=False, random_state=None, shrinking=True,
    tol=0.001, verbose=False)

```

0.907804705239104

AUC: 0.864

AUC: 0.894

AUC: 0.904

AUC: 0.908

AUC: 0.878

#####

AUC from train data #####

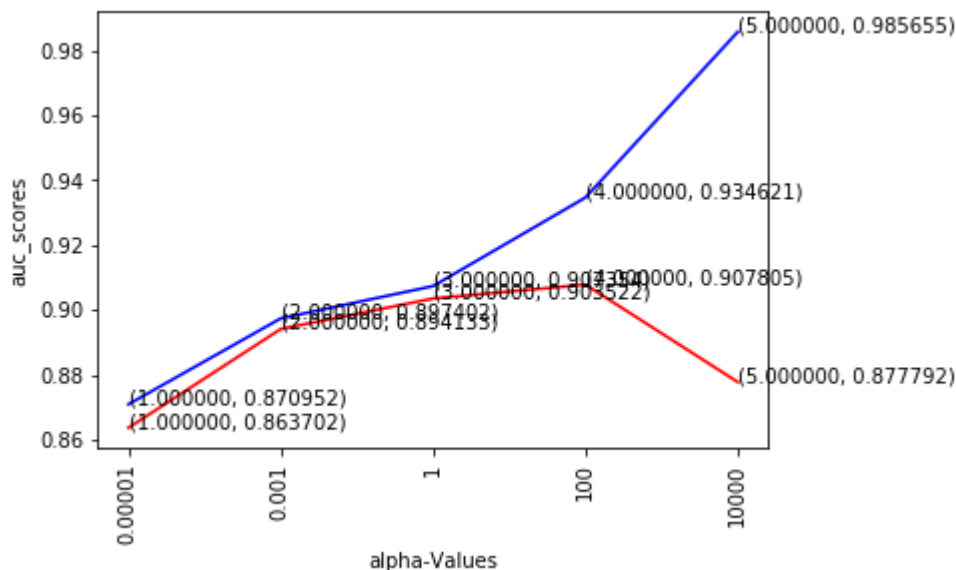
AUC: 0.871

AUC: 0.897

AUC: 0.907

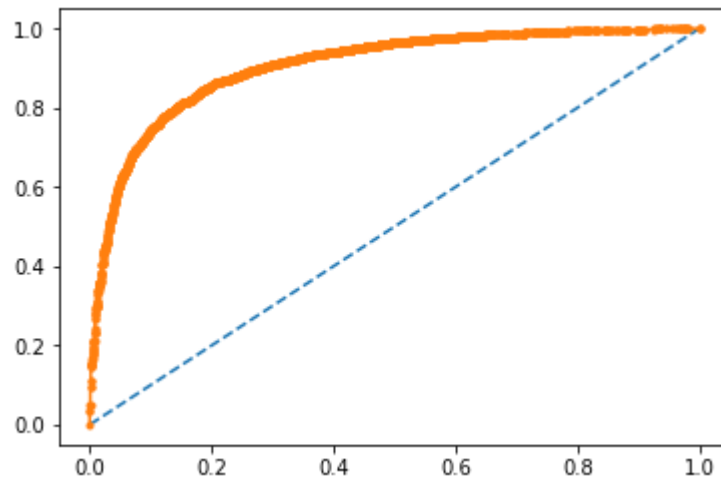
AUC: 0.935

AUC: 0.986



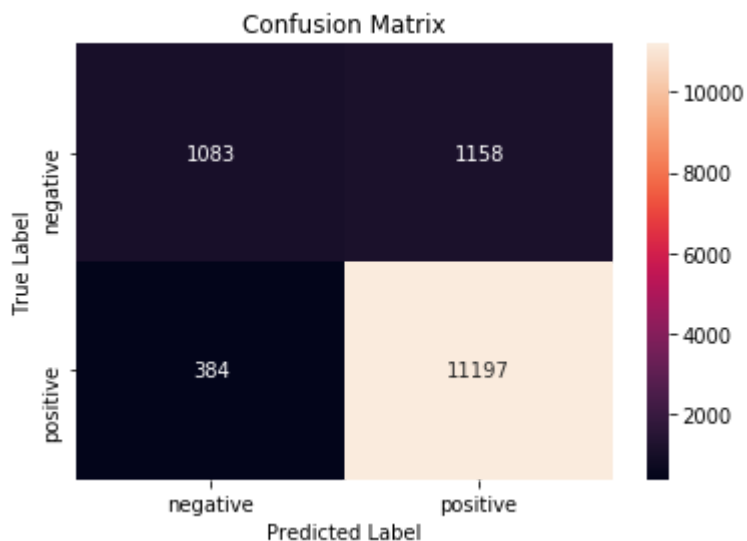
```
In [297]: 1 con_mat,clf = svc_results(100,sent_vectors,test_sent_vectors,y_1,y_test)
```

The accuracy of the SVC classifier for alpha = 100.000000 is 88.843872%
AUC: 0.906



Observation: Model predicted with accuracy 88 % with AUC: 0.906

```
In [298]: 1 showHeatMap(con_mat)
```



Observation: My model predicted 384 + 1158 points wrongly

[5.2.4] Applying RBF SVM on TFIDF W2V, SET 4

```
In [299]: 1 X_train, X_test, y_1, y_test = cross_validation.train_test_split(preprocessed
```

```
In [300]: 1 model = TfidfVectorizer()
2 X_train_transformed = model.fit_transform(X_train)
3
4 dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
```

```
In [301]: 1 # Train your own Word2Vec model using your own text corpus
2 i=0
3 list_of_sentence=[]
4 for sentence in X_train:
5     list_of_sentence.append(sentence.split())
```

```
In [302]: 1 w2v_model=Word2Vec(list_of_sentence,min_count=5,size=50, workers=4)
2 w2v_words = list(w2v_model.wv.vocab)
```

```
In [303]: 1 # TF-IDF weighted Word2Vec
2 tfidf_feat = model.get_feature_names() # tfidf words/col-names
3 # final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val
4
5 tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is stored i
6 row=0;
7 for sent in tqdm(list_of_sentence): # for each review/sentence
8     sent_vec = np.zeros(50) # as word vectors are of zero length
9     weight_sum =0; # num of words with a valid vector in the sentence/review
10    for word in sent: # for each word in a review/sentence
11        if word in w2v_words and word in tfidf_feat:
12            vec = w2v_model.wv[word]
13            # tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
14            # to reduce the computation we are
15            # dictionary[word] = idf value of word in whole courpus
16            # sent.count(word) = tf valeus of word in this review
17            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
18            sent_vec += (vec * tf_idf)
19            weight_sum += tf_idf
20    if weight_sum != 0:
21        sent_vec /= weight_sum
22    tfidf_sent_vectors.append(sent_vec)
23    row += 1
```

100%|██████████| 32249/32249 [19:22<00:00, 27.74it/s]

```
In [304]: 1 i=0
2 list_of_test_sentence=[]
3 for sentence in X_test:
4     list_of_test_sentence.append(sentence.split())
```

```

In [305]: 1 # TF-IDF weighted Word2Vec
2 tfidf_feat = model.get_feature_names() # tfidf words/col-names
3 # final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val
4
5 tfidf_test_sent_vectors = []; # the tfidf-w2v for each sentence/review is sto
6 row=0;
7 for sent in tqdm(list_of_test_sentence): # for each review/sentence
8     sent_vec = np.zeros(50) # as word vectors are of zero length
9     weight_sum =0; # num of words with a valid vector in the sentence/review
10    for word in sent: # for each word in a review/sentence
11        if word in w2v_words and word in tfidf_feat:
12            vec = w2v_model.wv[word]
13            # tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
14            # to reduce the computation we are
15            # dictionary[word] = idf value of word in whole corpus
16            # sent.count(word) = tf value of word in this review
17            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
18            sent_vec += (vec * tf_idf)
19            weight_sum += tf_idf
20    if weight_sum != 0:
21        sent_vec /= weight_sum
22    tfidf_test_sent_vectors.append(sent_vec)
23    row += 1

```

100%|██████████| 13822/13822 [07:58<00:00, 26.55it/s]

```

In [306]: 1 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_train, y_1, test
2
3 i=0
4 list_of_cv_sentence=[]
5 for sentence in X_tr:
6     list_of_cv_sentence.append(sentence.split())
7
8
9 tfidf_feat = model.get_feature_names() # tfidf words/col-names
10 # final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val
11
12 tfidf_cv_sent_vectors = []; # the tfidf-w2v for each sentence/review is store
13 row=0;
14 for sent in tqdm(list_of_cv_sentence): # for each review/sentence
15     sent_vec = np.zeros(50) # as word vectors are of zero length
16     weight_sum =0; # num of words with a valid vector in the sentence/review
17     for word in sent: # for each word in a review/sentence
18         if word in w2v_words and word in tfidf_feat:
19             vec = w2v_model.wv[word]
20             # tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
21             # to reduce the computation we are
22             # dictionary[word] = idf value of word in whole courpus
23             # sent.count(word) = tf valeus of word in this review
24             tf_idf = dictionary[word]*(sent.count(word)/len(sent))
25             sent_vec += (vec * tf_idf)
26             weight_sum += tf_idf
27         if weight_sum != 0:
28             sent_vec /= weight_sum
29         tfidf_cv_sent_vectors.append(sent_vec)
30         row += 1
31
32 i=0
33 list_of_cv_test_sentence=[]
34 for sentence in X_cv:
35     list_of_cv_test_sentence.append(sentence.split())
36
37
38 tfidf_cv_test_sent_vectors = []; # the tfidf-w2v for each sentence/review is
39 row=0;
40 for sent in tqdm(list_of_cv_test_sentence): # for each review/sentence
41     sent_vec = np.zeros(50) # as word vectors are of zero length
42     weight_sum =0; # num of words with a valid vector in the sentence/review
43     for word in sent: # for each word in a review/sentence
44         if word in w2v_words and word in tfidf_feat:
45             vec = w2v_model.wv[word]
46             # tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
47             # to reduce the computation we are
48             # dictionary[word] = idf value of word in whole courpus
49             # sent.count(word) = tf valeus of word in this review
50             tf_idf = dictionary[word]*(sent.count(word)/len(sent))
51             sent_vec += (vec * tf_idf)
52             weight_sum += tf_idf
53         if weight_sum != 0:
54             sent_vec /= weight_sum
55         tfidf_cv_test_sent_vectors.append(sent_vec)
56         row += 1

```

```

57
58
59 tuned_parameters = [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
60
61 #Using GridSearchCV
62 model = GridSearchCV(SVC(), tuned_parameters, scoring = 'roc_auc', cv=5)
63 model.fit(tfidf_cv_sent_vectors, y_tr)
64
65 print(model.best_estimator_)
66 print(model.score(tfidf_cv_test_sent_vectors, y_cv))
67
68 check_trade_off(tfidf_cv_sent_vectors,tfidf_cv_test_sent_vectors,y_tr,y_cv)
69

```

```

100%|██████████| 22574/22574 [10:36<00:00, 35.47it/s]

```

```

100%|██████████| 9675/9675 [06:14<00:00, 17.69it/s]

```

```

SVC(C=100, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto', kernel='rbf',
    max_iter=-1, probability=False, random_state=None, shrinking=True,
    tol=0.001, verbose=False)

```

```
0.8963067589437549
```

```
AUC: 0.814
```

```
AUC: 0.875
```

```
AUC: 0.884
```

```
AUC: 0.896
```

```
AUC: 0.856
```

```
#####
```

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

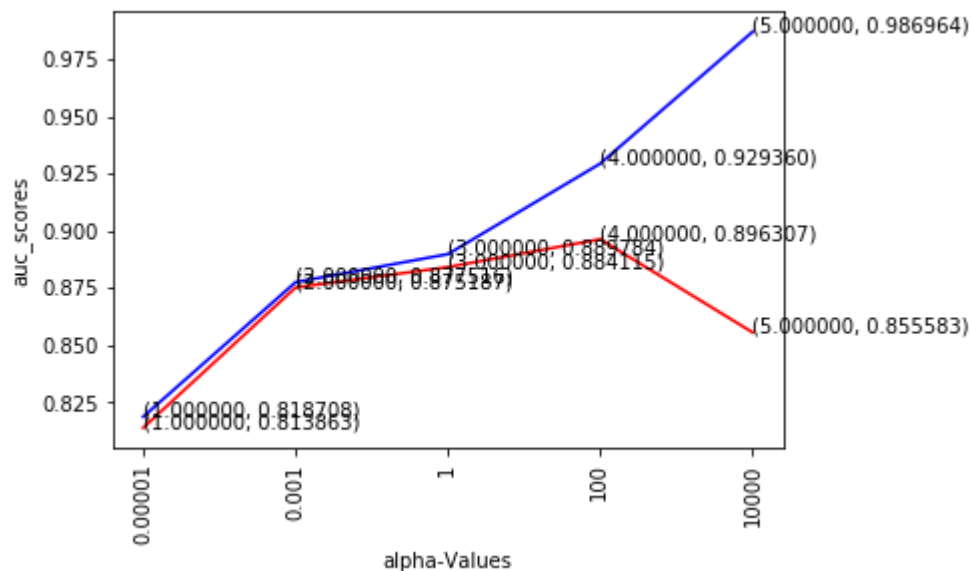
```
AUC: 0.819
```

```
AUC: 0.878
```

```
AUC: 0.890
```

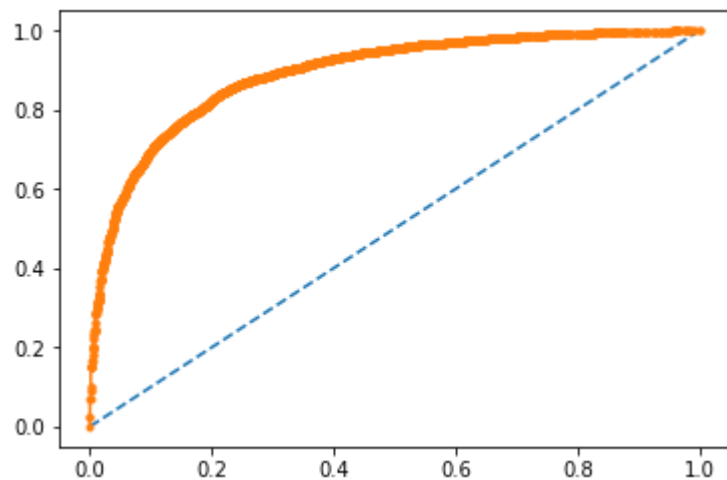
```
AUC: 0.929
```

```
AUC: 0.987
```

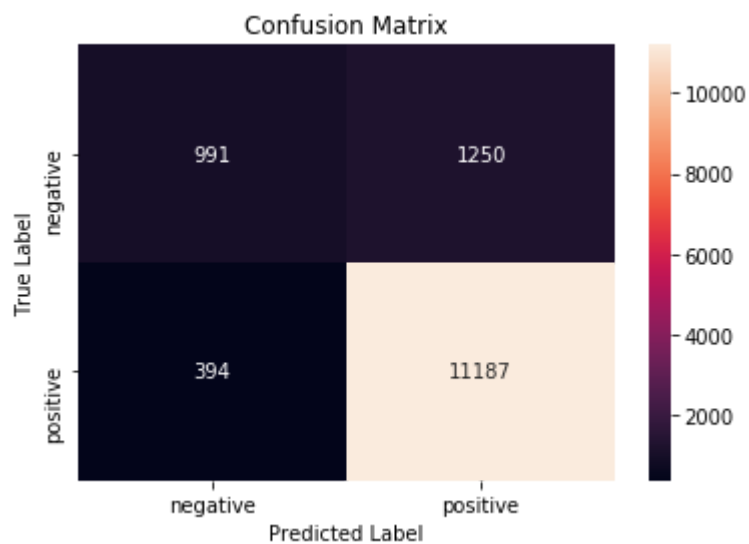



```
In [307]: 1 con_mat,clf = svc_results(100,tfidf_sent_vectors,tfidf_test_sent_vectors,y_1,
```

The accuracy of the SVC classifier for alpha = 100.000000 is 88.105918%
AUC: 0.892



```
In [308]: 1 showHeatMap(con_mat)
```



Repeat with extra features

```
In [38]: 1 from sklearn.cross_validation import train_test_split
2 from sklearn.svm import SVC
3 from sklearn.metrics import accuracy_score
4 from sklearn.cross_validation import cross_val_score
5 from collections import Counter
6 from sklearn.metrics import accuracy_score
7 from sklearn import cross_validation
8 from sklearn.grid_search import GridSearchCV
9 from sklearn.calibration import CalibratedClassifierCV
10 from sklearn.linear_model import SGDClassifier
11 import warnings
12 warnings.filterwarnings("ignore")
```

C:\Users\sujpanda\Anaconda3\lib\site-packages\sklearn\cross_validation.py:41: DeprecationWarning: This module was deprecated in version 0.18 in favor of the model_selection module into which all the refactored classes and functions are moved. Also note that the interface of the new CV iterators are different from that of this module. This module will be removed in 0.20.

"This module will be removed in 0.20.", DeprecationWarning)

C:\Users\sujpanda\Anaconda3\lib\site-packages\sklearn\grid_search.py:42: DeprecationWarning: This module was deprecated in version 0.18 in favor of the model_selection module into which all the refactored classes and functions are moved. This module will be removed in 0.20.

DeprecationWarning)

```
In [39]: 1 mylen = np.vectorize(len)
2 newarr = mylen(preprocessed_summary)
```

```
In [40]: 1 newproce_reviews = np.asarray(preprocessed_reviews)
```

```
In [41]: 1 newproce_summary = np.asarray(preprocessed_summary)
```

```
In [42]: 1 df = pd.DataFrame({'desc':newproce_reviews, 'summary':newproce_summary, 'len':
```

```
In [43]: 1 df.head()
```

Out[43]:

	desc	summary	len
0	dogs loves chicken product china wont buying a...	made china	10
1	dogs love saw pet store tag attached regarding...	dog lover delites	17
2	product available victor traps unreal course t...	thirty bucks	12
3	used victor fly bait seasons ca not beat great...	flies begone	12
4	received shipment could hardly wait try produc...	wow make islickers	18

```
In [44]: 1 X_1, X_test, y_1, y_test = cross_validation.train_test_split(df, final['Score
```

In [45]:

```

1 import scipy
2 count_vect = CountVectorizer()
3 final_counts = count_vect.fit_transform(X_1['desc'])
4 final_test_count = count_vect.transform(X_test['desc'])
5
6 # split the train data set into cross validation train and cross validation test
7 X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_1, y_1, test_size=0.2)
8
9 final_counts_tr_cv = count_vect.transform(X_tr['desc'])
10 final_test_count_cv = count_vect.transform(X_cv['desc'])
11
12 from scipy.sparse import csr_matrix, issparse
13
14 #####Adding len as feature#####
15 #if issparse(final_counts_tr_cv):
16     #print('sparse matrix')
17 len_sparse = scipy.sparse.coo_matrix(X_tr['len'])
18 len_sparse = len_sparse.transpose()
19
20 final_counts_tr_cv = scipy.sparse.hstack([final_counts_tr_cv, len_sparse])
21 print(final_counts_tr_cv.shape)
22
23 len_test_sparse = scipy.sparse.coo_matrix(X_cv['len'])
24 len_test_sparse = len_test_sparse.transpose()
25 final_test_count_cv = scipy.sparse.hstack([final_test_count_cv, len_test_sparse])
26 print("final_counts_tr_cv.shape after length = ", final_counts_tr_cv.shape)
27
28 #####Adding summary as feature#####
29 final_summary_count = count_vect.transform(X_tr['summary'])
30 final_test_summary_count_cv = count_vect.transform(X_cv['summary'])
31 columns=count_vect.get_feature_names()
32
33 print("sujet", final_summary_count[:,12].shape)
34 final_counts_tr_cv = scipy.sparse.hstack([final_counts_tr_cv, final_summary_count])
35 print("final_counts_tr_cv.shape after f1= ", final_counts_tr_cv.shape)
36
37 final_test_count_cv = scipy.sparse.hstack([final_test_count_cv, final_test_summary_count_cv])
38
39
40 final_counts_tr_cv = scipy.sparse.hstack([final_counts_tr_cv, final_summary_count])
41 print("final_counts_tr_cv.shape after f2= ", final_counts_tr_cv.shape)
42
43
44 final_test_count_cv = scipy.sparse.hstack([final_test_count_cv, final_test_summary_count_cv])
45
46 #####finding the new C #####
47
48 tuned_parameters = [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
49
50 #Using GridSearchCV
51 model = GridSearchCV(SVC(), tuned_parameters, scoring = 'roc_auc', cv=5)
52 model.fit(final_counts_tr_cv, y_tr)
53
54 print(model.best_estimator_)
55 print(model.score(final_test_count_cv, y_cv))
56

```

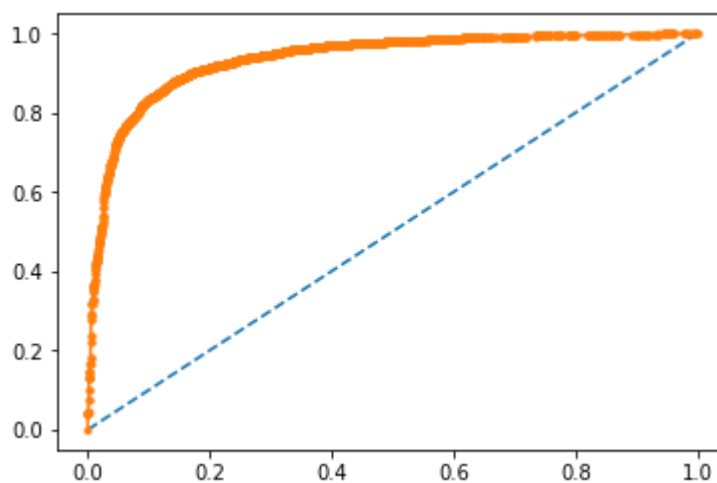
```

(22574, 33291)
final_counts_tr_cv.shape after length = (22574, 33291)
sujet (22574, 1)
final_counts_tr_cv.shape after f1= (22574, 33292)
final_counts_tr_cv.shape after f2= (22574, 33293)
SVC(C=100, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto', kernel='rbf',
    max_iter=-1, probability=False, random_state=None, shrinking=True,
    tol=0.001, verbose=False)
0.9329803897735287

```

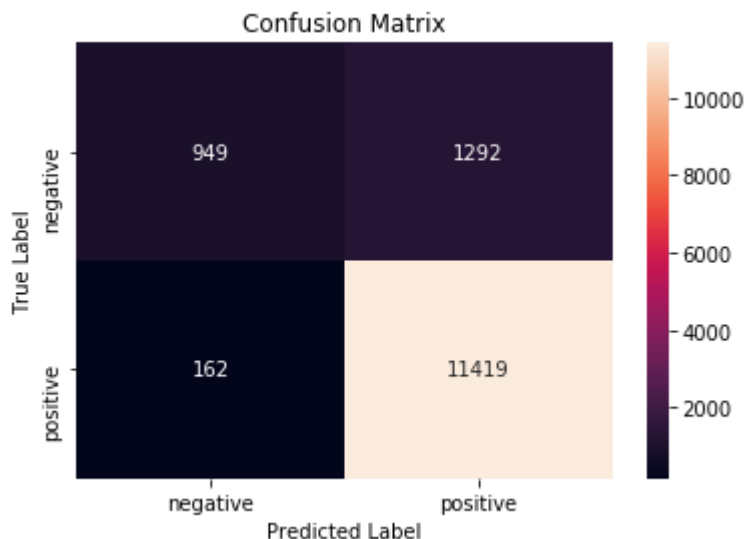
```
In [46]: 1 con_mat,clf = svc_results(100,final_counts,final_test_count,y_1,y_test)
```

The accuracy of the SVC classifier for alpha = 100.000000 is 89.480538%
AUC: 0.936



Observation: Model predicted with accuracy 89% with AUC: 936

```
In [47]: 1 showHeatMap(con_mat)
```



Observation: MY model predicted 162 + 1292 points wrongly

[6] Conclusions

Method	No of samples	Algorithm	alpha or C value	accuray	AUC Score	regularizer
BOW	50000	SGD	0.0001	90	0.909	l1
BOW	50000	SGD	0.0001	89	0.922	l2
TF-IDF	50000	SGD	0.0001	90	0.938	l1
TF-IDF	50000	SGD	0.0001	91	0.957	l2
AVG W2VEC	50000	SGD	0.0001	87	0.888	l1
AVG W2VEC	50000	SGD	0.0001	87	0.880	l2
TF_IDF AVG W2VEC	50000	SGD	0.0001	83	0.840	l1
TF_IDF AVG W2VEC	50000	SGD	0.0001	86	0.852	l2
BOW1	50000	SGD	100	88	0.906	
BOW	50000	SGD	10000	88	0.907	
AVG W2VEC	50000	SGD	100	88	0.906	
TF_IDF AVG W2VEC	50000	SGD	0.0001	89	0.936	

In []:

1