

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 [3]:

```
%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\Hi\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 [4]:

```

# using SQLite Table to read data.
con = sqlite3.connect('C:/ML/amazon-fine-food-reviews/database.sqlite')

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

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

filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3 LIMIT 100000""", con)

# Give reviews with Score>3 a positive rating(1), and reviews with a score<3 a negative rating(0).
def partition(x):
    if x < 3:
        return 0
    return 1

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

```

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

Out[4]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	1	1303862400	Good Quality Dog Food
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	0	1346976000	Not as Advertised

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary
2	3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1	1	1219017600	"Delight" says it all

In [5]:

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

In [6]:

```
print(display.shape)
display.head()
```

(80668, 7)

Out [6]:

	UserId	ProductId	ProfileName	Time	Score	Text	COUNT(*)
0	#oc-R115TNMSPFT9I7	B005ZBZLT4	Breyton	1331510400	2	Overall its just OK when considering the price...	2
1	#oc-R11D9D7SHXIJB9	B005HG9ESG	Louis E. Emory "hoppy"	1342396800	5	My wife has recurring extreme muscle spasms, u...	3
2	#oc-R11DNU2NBKQ23Z	B005ZBZLT4	Kim Cieszykowski	1348531200	1	This coffee is horrible and unfortunately not ...	2
3	#oc-R11O5J5ZVQE25C	B005HG9ESG	Penguin Chick	1346889600	5	This will be the bottle that you grab from the...	3
4	#oc-R12KPBODL2B5ZD	B007OSBEV0	Christopher P. Presta	1348617600	1	I didnt like this coffee. Instead of telling y...	2

In [7]:

```
display[display['UserId']=='AZY10LLTJ71NX']
```

Out [7]:

	UserId	ProductId	ProfileName	Time	Score	Text	COUNT(*)
80638	AZY10LLTJ71NX	B001ATMQK2	undertheshrine "undertheshrine"	1296691200	5	I bought this 6 pack because for the price tha...	5

In [8]:

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

Out [8]:

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 [9]:

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

Out[9]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summ
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA ^T VANII WAFE
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA ^T VANII WAFE
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA ^T VANII WAFE
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA ^T VANII WAFE
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA ^T VANII WAFE

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 [10]:

```
#Sorting data according to ProductId in ascending order
sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True, inplace=False, kind='quicksort', na_position='last')
```

In [11]:

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

Out[11]:

(87775, 10)

In [12]:

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

Out[12]:

87.775

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

In [13]:

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

display.head()
```

Out[13]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary
0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	1	5	1224892800	Bought This for My Son at College
1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	2	4	1212883200	Pure cocoa taste with crunchy almonds inside

In [14]:

```
final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

In [15]:

```
#Before starting the next phase of preprocessing lets see the number of entries left
print(final.shape)

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

(87773, 10)

Out[15]:

```
1    73592
0    14181
Name: Score, dtype: int64
```

[3] Preprocessing

[3.1]. Preprocessing Review Text

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

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

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

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

In [16]:

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

sent_1000 = final['Text'].values[1000]
print(sent_1000)
print("="*50)

sent_1500 = final['Text'].values[1500]
print(sent_1500)
print("="*50)

sent_4900 = final['Text'].values[4900]
print(sent_4900)
print("="*50)
```

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

=====

The Candy Blocks were a nice visual for the Lego Birthday party but the candy has little taste to it. Very little of the 2 lbs that I bought were eaten and I threw the rest away. I would not buy the candy again.

=====

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

=====

My dog LOVES these treats. They tend to have a very strong fish oil smell. So if you are afraid of the fishy smell, don't get it. But I think my dog likes it because of the smell. These treats are really small in size. They are great for training. You can give your dog several of these without worrying about him over eating. Amazon's price was much more reasonable than any other retailer. You can buy a 1 pound bag on Amazon for almost the same price as a 6 ounce bag at other retailers. It's definitely worth it to buy a big bag if your dog eats them a lot.

=====

In [17]:

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

print(sent_0)
```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its ver y hard to find any chicken products made in the USA but they are out there, but this one isnt. It s 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 [18]:

```
# https://stackoverflow.com/questions/16206380/python-beautifulsoup-how-to-remove-all-tags-from-an-element
from bs4 import BeautifulSoup

soup = BeautifulSoup(sent_0, 'lxml')
```

```

text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_1000, 'lxml')
text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_1500, 'lxml')
text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_4900, 'lxml')
text = soup.get_text()
print(text)

```

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

=====

The Candy Blocks were a nice visual for the Lego Birthday party but the candy has little taste to it. Very little of the 2 lbs that I bought were eaten and I threw the rest away. I would not buy the candy again.

=====

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

=====

My dog LOVES these treats. They tend to have a very strong fish oil smell. So if you are afraid of the fishy smell, don't get it. But I think my dog likes it because of the smell. These treats are really small in size. They are great for training. You can give your dog several of these without worrying about him over eating. Amazon's price was much more reasonable than any other retailer. Y ou can buy a 1 pound bag on Amazon for almost the same price as a 6 ounce bag at other retailers. It's definitely worth it to buy a big bag if your dog eats them a lot.

In [19]:

```

# https://stackoverflow.com/a/47091490/4084039
import re

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

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

```

In [20]:

```

sent_1500 = decontracted(sent_1500)
print(sent_1500)
print("="*50)

```

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

=====

In [21]:

```

#remove words with numbers python: https://stackoverflow.com/a/18082370/4084039
sent_0 = re.sub(r"\S*\d\S*", "", sent_0).strip()
print(sent_0)

```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its ver y hard to find any chicken products made in the USA but they are out there, but this one isnt. It s 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 [22]:

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

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

In [23]:

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

stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "y
ou're", "you've",\
    "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his',
'himself', \
    'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them',
'their',\
    'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll",
'these', 'those', \
    'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having',
'do', 'does', \
    'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until',
while', 'of', \
    'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during',
'before', 'after',\
    'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under',
, 'again', 'further',\
    'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'e
ach', 'few', 'more',\
    'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
    's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll'
, 'm', 'o', 're', \
    've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "do
esn't", 'hadn',\
    "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn',
'mightn't", 'mustn',\
    "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn',
'wasn't", 'weren', "weren't", \
    'won', "won't", 'wouldn', "wouldn't"]])
```

In [24]:

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

[illegible]

In [25]:

```
preprocessed_reviews[1500]
```



```
Out[25]:
```

```
'way hot blood took bite jig lol'
```

[3.2] Preprocessing Review Summary

Similarly you can do preprocessing for review summary also.

[4] Featurization

[4.1] BAG OF WORDS

```
In [26]:
```

```
#BoW
count_vect = CountVectorizer(max_features = 2000,min_df = 10) #in scikit-learn
count_vect.fit(preprocessed_reviews)
print("some feature names ", count_vect.get_feature_names()[0:10])
print('='*50)

final_counts = count_vect.transform(preprocessed_reviews)
print("the type of count vectorizer ",type(final_counts))
print("the shape of out text BOW vectorizer ",final_counts.get_shape())
print("the number of unique words ", final_counts.get_shape()[1])
```

```
some feature names  ['able', 'absolute', 'absolutely', 'according', 'acid', 'acidic', 'across', 'actual', 'actually', 'add']
=====
```

```
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (87773, 2000)
the number of unique words  2000
```

[4.2] Bi-Grams and n-Grams.

```
In [27]:
```

```
#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 numebtrs 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])
```

```
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (87773, 5000)
the number of unique words including both unigrams and bigrams  5000
```

[4.3] TF-IDF

```
In [28]:
```

```
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10,max_features=5000)
tf_idf_vect.fit(preprocessed_reviews)
print("some sample features(unique words in the corpus)",tf_idf_vect.get_feature_names()[0:10])
print('='*50)
```

```
print('='*50)

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

some sample features(unique words in the corpus) ['ability', 'able', 'able buy', 'able find', 'able get', 'absolute', 'absolute favorite', 'absolutely', 'absolutely best', 'absolutely delicious']

```
=====
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer (87773, 5000)
the number of unique words including both unigrams and bigrams 5000
```

[4.4] Word2Vec

In [29]:

```
# Train your own Word2Vec model using your own text corpus
i=0
list_of_sentence=[]
for sentence in preprocessed_reviews:
    list_of_sentence.append(sentence.split())
```

In [30]:

```
# Using Google News Word2Vectors

# in this project we are using a pretrained model by google
# its 3.3G file, once you load this into your memory
# it occupies ~9Gb, so please do this step only if you have >12G of ram
# we will provide a pickle file wich contains a dict ,
# and it contains all our courpus words as keys and model[word] as values
# To use this code-snippet, download "GoogleNews-vectors-negative300.bin"
# from https://drive.google.com/file/d/0B7XkCwpI5KDYNlNUTTlSS21pQmM/edit
# it's 1.9GB in size.

# http://kavita-ganesan.com/gensim-word2vec-tutorial-starter-code/#.W17SRFAzZPY
# you can comment this whole cell
# or change these variable according to your need

want_to_train_w2v = True

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

```
[('fantastic', 0.8491621017456055), ('awesome', 0.8436824679374695), ('good', 0.8203500509262085),
('excellent', 0.8102934956550598), ('terrific', 0.7967278957366943), ('wonderful',
0.7863171100616455), ('perfect', 0.7452594637870789), ('amazing', 0.7439837455749512), ('nice', 0.
7136707305908203), ('decent', 0.690740168094635)]
```

```
=====
[('greatest', 0.8035088777542114), ('tastiest', 0.7343665361404419), ('best', 0.719451904296875),
('disgusting', 0.6661298871040344), ('nastiest', 0.6288201808929443), ('horrible',
0.6175878643989563), ('closest', 0.6150859594345093), ('wins', 0.6050095558166504), ('awful', 0.60
21138429641724), ('terrible', 0.5984127521514893)]
```

In [31]:

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

number of words that occurred minimum 5 times 17386
sample words ['dogs', 'loves', 'chicken', 'product', 'china', 'wont', 'buying', 'anymore',

```
'hard', 'find', 'products', 'made', 'usa', 'one', 'isnt', 'bad', 'good', 'take', 'chances',  
'till', 'know', 'going', 'imports', 'love', 'saw', 'pet', 'store', 'tag', 'attached', 'regarding',  
'satisfied', 'safe', 'infestation', 'literally', 'everywhere', 'flying', 'around', 'kitchen',  
'bought', 'hoping', 'least', 'get', 'rid', 'weeks', 'fly', 'stuck', 'squishing', 'buggers', 'succe  
ss', 'rate']
```

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

[4.4.1.1] Avg W2v

In [32]:

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

100%|██| 87773/87773 [05:06<00:00, 286.49it/s]

87773
50

[4.4.1.2] TFIDF weighted W2v

In [33]:

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]  
model = TfidfVectorizer(min_df=10, max_features=5000)  
tf_idf_matrix = model.fit_transform(preprocessed_reviews)  
# we are converting a dictionary with word as a key, and the idf as a value  
dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
```

In [34]:

```
# TF-IDF weighted Word2Vec  
tfidf_feat = model.get_feature_names() # tfidf words/col-names  
# final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val = tfidf  
  
tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is stored in this list  
row=0;  
for sent in tqdm(list_of_sentence): # 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.append(sent_vec)
```

```
100%|███████████| 87773/87773 [10:23<00:00, 140.67it/s]
```

- [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 you train data, and apply the method `transform()` on cv/test data.
4. For more details please go through this [link](#).

Applying Logistic Regression

[5.1] Logistic Regression on BOW, SET 1

[5.1.1] Applying Logistic Regression with L1 regularization on BOW, SET 1

In [91]:

```
# ===== loading libraries =====
import pdb
import numpy as np
import pandas as pd
import seaborn as sns
import copy
import matplotlib.pyplot as plt
import sklearn
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
from sklearn.metrics import accuracy_score
from sklearn.metrics import roc_auc_score
from sklearn.metrics import confusion_matrix
from collections import Counter
import scikitplot.metrics as skplt
from sklearn.linear_model import LogisticRegression
# =====
```

5.1.12 Splitting the data converting to bag of words

In [92]:

```
#Splitting entire data to train,test and cross validation
X=np.array(preprocessed_reviews)
y = np.array(final['Score'])

## split the data set into train and test
X_1, X_test, y_1, y_test = train_test_split(X, y, test_size=0.3, random_state=1)

# split the train data set into cross validation train and cross validation test
X_tr, X_cv, y_tr, y_cv = train_test_split(X_1, y_1, test_size=0.3,random_state=1)

#converting Reviews to Bag of words after splitting to avoid data leakage problem
count_vect = CountVectorizer(max_features = 5000,min_df = 10)
final_X_tr=count_vect.fit_transform(X_tr)
final_X_test=count_vect.transform(X_test)
final_X_cv=count_vect.transform(X_cv)
```

In [98]:

```
#Calculating for finding Best C
#predic_proba reference:
#https://stackoverflow.com/questions/37089177/probability-prediction-method-of-
kneighborsclassifier-returns-only-0-and-1
#https://discuss.analyticsvidhya.com/t/what-is-the-difference-between-predict-and-predict_proba/67
376/3
roc_tr=[]
roc_cv=[]
max_auc_score=0
best_c=0
tuned_parameters =[10**-4, 10**-2, 10**0, 10**2, 10**4]
```

```

for i in tuned_parameters:
    clf=LogisticRegression(C=i, penalty='l1', class_weight='balanced')

    # fitting the model on train data
    clf.fit(final_X_tr,y_tr)

    # predict the response on the traininig
    pred_tr = clf.predict_proba(final_X_tr)
    pred_tr=(pred_tr)[:,1]
    roc_tr.append(roc_auc_score(y_tr,pred_tr))

    # predict the response on the crossvalidation
    pred_cv = clf.predict_proba(final_X_cv)
    pred_cv=(pred_cv)[:,1]
    roc_cv.append(roc_auc_score(y_cv,pred_cv))

    #finding best c using loop
    if roc_auc_score(y_cv,pred_cv)>max_auc_score:
        best_c=i
        max_auc_score=roc_auc_score(y_cv,pred_cv)

print(best_c)
print(max_auc_score)
c_llset1=best_c
auc_set1=max_auc_score

```

```

1
0.9247879668432126

```

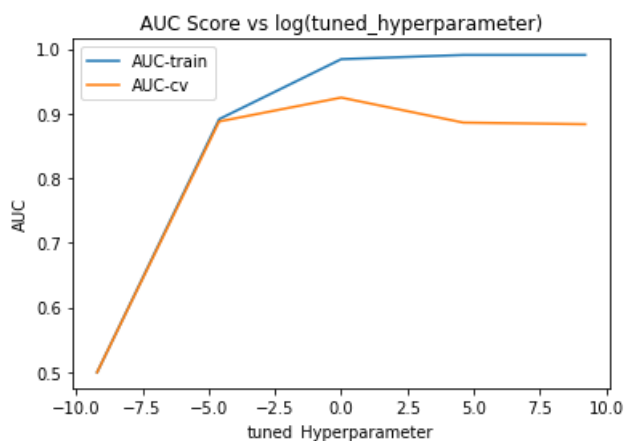
5.1.14 Curve plotting between AUC of cv and train with log(tuned_parameters)

In [99]:

```

# plotting curve between between AUC of cv and train with log of tuned parameter
logalpha=np.log(tuned_parameters)
plt.plot(logalpha,roc_tr,label="AUC-train")
plt.plot(logalpha,roc_cv ,label="AUC-cv")
plt.legend()
plt.xlabel('tuned_Hyperparameter')
plt.ylabel('AUC')
plt.title('AUC Score vs log(tuned_hyperparameter) ')
plt.show()

```



5.1.15 Training the model with the obtained best_c and plotting Roc curve

In [100]:

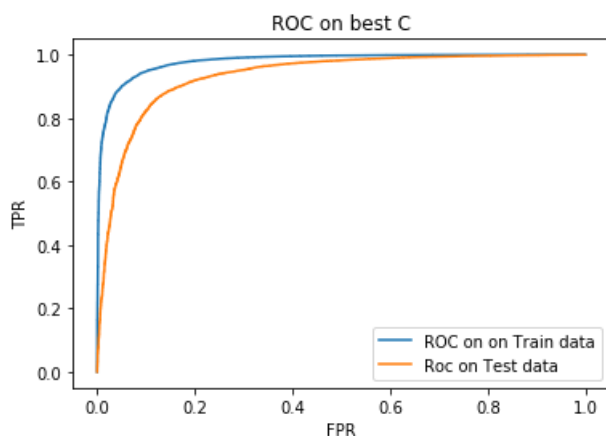
```

#1) Training the model using best C
clf=LogisticRegression(C=c_llset1, penalty='l1')
clf.fit(final_X_tr, y_tr)
#predicting probability on Test data
pred_test = clf.predict_proba(final_X_test)
pred_test=(pred_test)[:,1]
#predicting probablity of Training data
pred_tr = clf.predict_proba(final_X_tr)

```

```
pred_tr=(pred_tr)[:,-1]
```

```
#2)Plotting Roc Curve
#Reference for finding fpr an tpr :
#https://www.programcreek.com/python/example/81207/sklearn.metrics.roc_curve
fpr_tr, tpr_tr, threshold_train = metrics.roc_curve(y_tr, pred_tr)
fpr_test, tpr_test, threshold_test = metrics.roc_curve(y_test, pred_test)
plt.plot(fpr_tr,tpr_tr ,label="ROC on on Train data")
plt.plot(fpr_test,tpr_test ,label="Roc on Test data")
plt.legend()
plt.title('ROC on best C')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



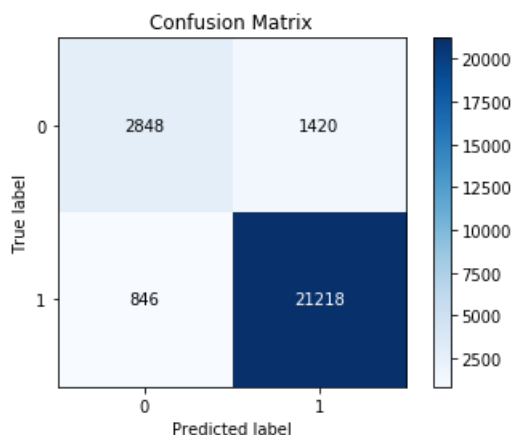
In [101]:

```
#plotting the confusion matrix
#Reference:
#https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html

prediction=clf.predict(final_X_test)
skplt.plot_confusion_matrix(y_test ,prediction)
```

Out[101]:

<matplotlib.axes._subplots.AxesSubplot at 0x9ea82b0b70>



[5.1.1.1] Calculating sparsity on weight vector obtained using L1 regularization on BOW, SET 1

In [102]:

```
lr = LogisticRegression(C=c_1lset1, penalty='l1');
lr.fit(final_X_tr, y_tr);
w = lr.coef_
print(np.count_nonzero(w))
```

[5.1.2] Applying Logistic Regression with L2 regularization on BOW, SET 1

In [103]:

```
#Calculating for finding Best C
#predic_proba reference:
#https://stackoverflow.com/questions/37089177/probability-prediction-method-of-
#neighborsclassifier-returns-only-0-and-1
#https://discuss.analyticsvidhya.com/t/what-is-the-difference-between-predict-and-predict_proba/67
#376/3
roc_tr=[]
roc_cv=[]
max_auc_score=0
best_c=0
tuned_parameters =[10**-4, 10**-2, 10**0, 10**2, 10**4]
for i in tuned_parameters:
    clf=LogisticRegression(C=i, penalty='l2', class_weight='balanced',max_iter=1000)
    # fitting the model on train data
    clf.fit(final_X_tr,y_tr)

    # predict the response on the traininig
    pred_tr = clf.predict_proba(final_X_tr)
    pred_tr=(pred_tr)[:,1]
    roc_tr.append(roc_auc_score(y_tr,pred_tr))

    # predict the response on the crossvalidation
    pred_cv = clf.predict_proba(final_X_cv)
    pred_cv=(pred_cv)[:,1]
    roc_cv.append(roc_auc_score(y_cv,pred_cv))

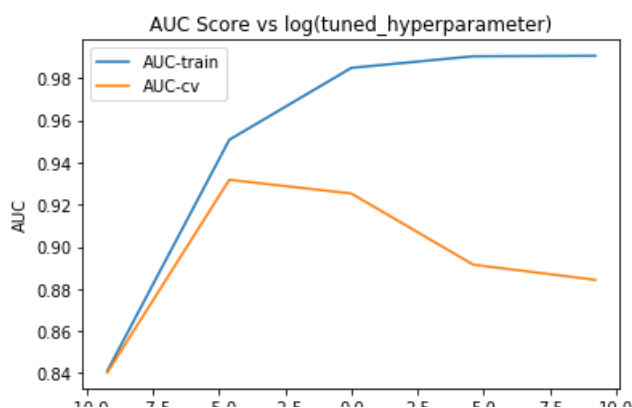
    #finding best c using loop
    if roc_auc_score(y_cv,pred_cv)>max_auc_score:
        best_c=i
        max_auc_score=roc_auc_score(y_cv,pred_cv)

print(best_c)
print(max_auc_score)
c_l2set1=best_c
auc_set2=max_auc_score
```

0.01
0.9318337422621068

In [104]:

```
# plotting curve between between AUC of cv and train with log of tuned parameter
logalpha=np.log(tuned_parameters)
plt.plot(logalpha,roc_tr,label="AUC-train")
plt.plot(logalpha,roc_cv ,label="AUC-cv")
plt.legend()
plt.xlabel('tuned_Hyperparameter')
plt.ylabel('AUC')
plt.title('AUC Score vs log(tuned_hyperparameter) ')
plt.show()
```

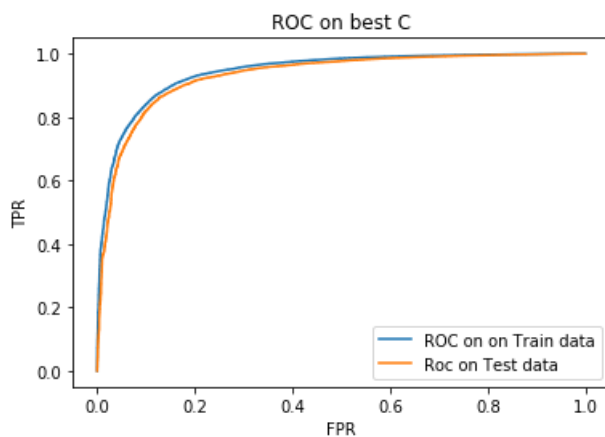


-10.0 -7.5 -5.0 -2.5 0.0 2.5 5.0 7.5 10.0
tuned_Hyperparameter

In [168]:

```
#1) Training the model using best C
clf=LogisticRegression(C=c_l2set1, penalty='l2')
clf.fit(final_X_tr, y_tr)
#predicting probability on Test data
pred_test = clf.predict_proba(final_X_test)
pred_test=(pred_test)[:,1]
#predicting probability of Training data
pred_tr = clf.predict_proba(final_X_tr)
pred_tr=(pred_tr)[:,1]

#2) Plotting Roc Curve
#Reference for finding fpr an tpr :
#https://www.programcreek.com/python/example/81207/sklearn.metrics.roc_curve
fpr_tr, tpr_tr, threshold_train = metrics.roc_curve(y_tr, pred_tr)
fpr_test, tpr_test, threshold_test = metrics.roc_curve(y_test, pred_test)
plt.plot(fpr_tr,tpr_tr ,label="ROC on on Train data")
plt.plot(fpr_test,tpr_test ,label="Roc on Test data")
plt.legend()
plt.title('ROC on best C')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



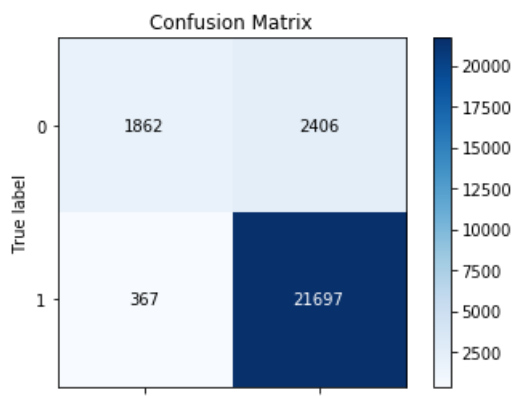
In [112]:

```
#plotting the confusion matrix
#Reference:
#https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html

prediction=clf.predict(final_X_test)
skplt.plot_confusion_matrix(y_test ,prediction)
```

Out[112]:

<matplotlib.axes._subplots.AxesSubplot at 0x9e92085898>



0 1
Predicted label

[5.1.2.1] Performing perturbation test (multicollinearity check) on BOW, SET 1

In [240]:

```
#a) Getting the weights before adding noise
weight_beforeperb=clf.coef_[0]
print(weight_beforeperb)
```

```
[-3.52877405e-05  1.48885337e-01  2.06881206e-02 ... -1.00581669e-02
 8.78842800e-03  1.42280663e-02]
```

In [241]:

```
#b) Add a noise to the X ( $X' = X + e$ ) and get the new data set  $X'$  (if  $X$  is a sparse matrix,
 $X.data+=e$ )
noisedata_tr=copy.deepcopy(final_X_tr)
noisedata_tr.data = noisedata_tr.astype(float)
noisedata_tr.data.data +=0.001
```

In [242]:

```
# c) fitting the model again on data  $X'$  and get the weights  $W'$ 
lr=LogisticRegression(C=c_l2set1, penalty='l2')
lr.fit(noisedata_tr.data, y_tr)
weight_afterperb=lr.coef_[0]
print(weight_afterperb)
```

```
[-4.05087435e-05  1.48864582e-01  2.06825298e-02 ... -1.00669616e-02
 8.80163402e-03  1.42326778e-02]
```

In [243]:

```
# d) Adding the small eps value(to eliminate the divisible by zero error) to  $W$  and  $W'$  i.e  $W=W+10^{-6}$ 
and  $W' = W'+10^{-6}$ 
weight_beforeperb = weight_beforeperb + 10**-6
weight_afterperb = weight_afterperb + 10**-6
```

In [244]:

```
#e) finding the % change between  $W$  and  $W'$ , percentage_change_vector = (|  $(W-W')$  | /  $(W)$  |)*100
list=[]
for i in range(len(weight_beforeperb)):
    list.append(((weight_beforeperb[i]-weight_afterperb[i])/weight_beforeperb[i])*100)

#Got help from AAIC members on how to use this list in percentiles
new_vector=np.absolute(list)
sortednew_vector=sorted(np.absolute(list))[:-1]
```

In [245]:

```
#f) calculate the 0th, 10th, 20th, 30th, ...100th percentiles
for i in range(0,110,10):
    print("{} percentile value is {}".format(i,str(np.percentile(sortednew_vector,i))))
```

```
0 percentile value is 7.898877957987713e-06
10 percentile value is 0.010380410177967836
20 percentile value is 0.019914074525295693
30 percentile value is 0.02914416207573957
40 percentile value is 0.03938667309106768
50 percentile value is 0.04999679477898631
60 percentile value is 0.06258026199425963
70 percentile value is 0.08065432131741687
80 percentile value is 0.10902068017333039
90 percentile value is 0.1941430827169683
100 percentile value is 48.73405175240942
```

In [246]:

```
#there is sudden rise in percentile from 90 to 100
#calculating percentile from 90 to 100
for i in range(90,101):
    print("{} percentile value is {}".format(i, str(np.percentile(sortednew_vector,i))))
```

```
90 percentile value is 0.1941430827169683
91 percentile value is 0.21340209076793107
92 percentile value is 0.23807063982860568
93 percentile value is 0.27191055015343313
94 percentile value is 0.32358577431264035
95 percentile value is 0.37353077314823874
96 percentile value is 0.46511878968660925
97 percentile value is 0.6413391216641606
98 percentile value is 1.0141104164840924
99 percentile value is 1.96321025044414
100 percentile value is 48.73405175240942
```

In [247]:

```
#there is sudden rise in percentile from 99 to 100
#calculating percentile from 99.1 to 100
for i in np.arange(0.0, 1.1, 0.1):
    print("{} percentile value is {}".format(99+i, np.percentile(sortednew_vector,99+i)))
```

```
99.0 percentile value is 1.96321025044414
99.1 percentile value is 2.204167275690913
99.2 percentile value is 2.4077504804947396
99.3 percentile value is 2.708441518653381
99.4 percentile value is 3.6194752893958846
99.5 percentile value is 4.681819243534795
99.6 percentile value is 6.644117839893287
99.7 percentile value is 7.890286294657444
99.8 percentile value is 12.403023555858185
99.9 percentile value is 15.22990656911503
100.0 percentile value is 48.73405175240942
```

In [248]:

```
#Got help from the AAIC members for how to print the feature names
#g)print the feature names whose % change is more than a threshold x(in our example it's 2.5)
all_features = count_vect.get_feature_names()
new_vector=new_vector.tolist()
index=new_vector.index(np.percentile(per_vector,100))
print(all_features[index])
```

trap

[5.1.3] Feature Importance on BOW, SET 1

[5.1.3.1] Top 10 important features of positive and negative class from SET 1

In [51]:

```
#Reference for top features is from statistics of machine learning by pratap dangeti:
#https://books.google.co.in/books?id=C-dDDwAAQBAJ&pg=PA216&lpg=PA216&dq=coefs_with_fns%5B:-(n+%2B+
1):-
1&source=bl&ots=j18t1prZXo&sig=ACfU3U2yz8v4v3SOiSrT7fBpJauJKe80DQ&hl=en&sa=X&ved=2ahUKewjzwK69x6DjI
XMBHYHqADIQ6AEwAXoECAGQAQ#v=onepage&q=coefs_with_fns%5B%3A-(n%20%2B%201)%3A-1&f=false
n=10
feature_names = count_vect.get_feature_names()
coefs_with_fns = sorted(zip(clf.coef_[0], feature_names))
top_n_coefs = zip(coefs_with_fns[:n], coefs_with_fns[-(n + 1):-1])
print("\tNegative-feature\t\t\t\t\tPositive-feature")
print("-----")
print("-----")
for (coef 1, fn 1), (coef 2, fn 2) in top_n_coefs:
```

```
print("\t%.4f\t%-15s\t\t\t\t%.4f\t%-15s" % (coef_1, fn_1, coef_2, fn_2))
```

Negative-feature	Positive-feature
-0.7374 disappointed	0.8834 great
-0.6397 worst	0.8561 delicious
-0.6152 money	0.8032 best
-0.6038 horrible	0.7086 perfect
-0.5570 awful	0.6938 loves
-0.5569 disappointing	0.6209 excellent
-0.5456 terrible	0.5776 wonderful
-0.5443 unfortunately	0.5389 favorite
-0.5039 waste	0.5385 love
-0.4915 stale	0.5371 nice

[5.2] Logistic Regression on TFIDF, SET 2

[5.2.1] Applying Logistic Regression with L1 regularization on TFIDF, SET 2

In [52]:

```
#Splitting entire data to train,test and cross validation
X=np.array(preprocessed_reviews)
y = np.array(final['Score'])

## split the data set into train and test
X_1, X_test, y_1, y_test = train_test_split(X, y, test_size=0.3, random_state=1)

# split the train data set into cross validation train and cross validation test
X_tr, X_cv, y_tr, y_cv = train_test_split(X_1, y_1, test_size=0.3, random_state=1)

#converting Reviews to tf_idf_vec
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2),min_df=10,max_features=5000)
final_X_tr=tf_idf_vect.fit_transform(X_tr)
final_X_test=tf_idf_vect.transform(X_test)
final_X_cv=tf_idf_vect.transform(X_cv)
```

In [53]:

```
#Calculating for finding Best C
#predic_proba reference:
#https://stackoverflow.com/questions/37089177/probability-prediction-method-of-kneighborsclassifier-returns-only-0-and-1
#https://discuss.analyticsvidhya.com/t/what-is-the-difference-between-predict-and-predict_proba/67376/3
roc_tr=[]
roc_cv=[]
max_auc_score=0
best_c=0
tuned_parameters=[10**-4, 10**-2, 10**0, 10**2, 10**4]
for i in tuned_parameters:
    clf=LogisticRegression(C=i, penalty='l1', class_weight='balanced')
    # fitting the model on train data
    clf.fit(final_X_tr,y_tr)

    # predict the response on the traininig
    pred_tr = clf.predict_proba(final_X_tr)
    pred_tr=(pred_tr)[:,-1]
    roc_tr.append(roc_auc_score(y_tr,pred_tr))

    # predict the response on the crossvalidation
    pred_cv = clf.predict_proba(final_X_cv)
    pred_cv=(pred_cv)[:,-1]
    roc_cv.append(roc_auc_score(y_cv,pred_cv))

    #finding best c using loop
    if roc_auc_score(y_cv,pred_cv)>max_auc_score:
        best_c=i
        max_auc_score=roc_auc_score(y_cv,pred_cv)

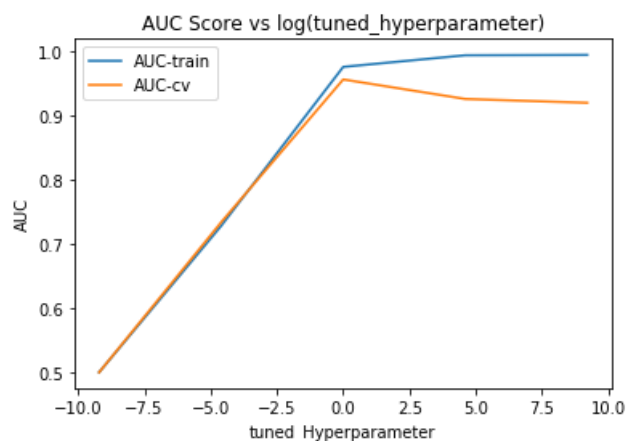
print(best_c)
```

```
print(max_auc_score)
c_l1set2=best_c
auc_set3=max_auc_score
```

```
1
0.9558029952014784
```

In [54]:

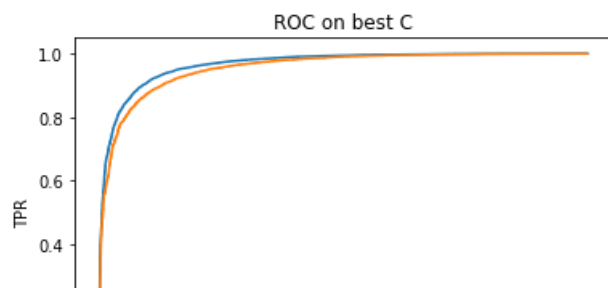
```
# plotting curve between between AUC of cv and train with log of tuned parameter
logalpha=np.log(tuned_parameters)
plt.plot(logalpha,roc_tr,label="AUC-train")
plt.plot(logalpha,roc_cv ,label="AUC-cv")
plt.legend()
plt.xlabel('tuned_Hyperparameter')
plt.ylabel('AUC')
plt.title('AUC Score vs log(tuned_hyperparameter)')
plt.show()
```

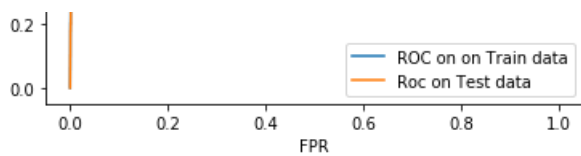


In [55]:

```
#1)Training the model using best C
clf=LogisticRegression(C=c_l1set2, penalty='l1')
clf.fit(final_X_tr, y_tr)
#predicting probability on Test data
pred_test = clf.predict_proba(final_X_test)
pred_test=(pred_test)[:,1]
#predicting probability of Training data
pred_tr = clf.predict_proba(final_X_tr)
pred_tr=(pred_tr)[:,1]

#2)Plotting Roc Curve
#Reference for finding fpr an tpr :
#https://www.programcreek.com/python/example/81207/sklearn.metrics.roc_curve
fpr_tr, tpr_tr, threshold_train = metrics.roc_curve(y_tr, pred_tr)
fpr_test, tpr_test, threshold_test = metrics.roc_curve(y_test, pred_test)
plt.plot(fpr_tr,tpr_tr ,label="ROC on on Train data")
plt.plot(fpr_test,tpr_test ,label="Roc on Test data")
plt.legend()
plt.title('ROC on best C')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```





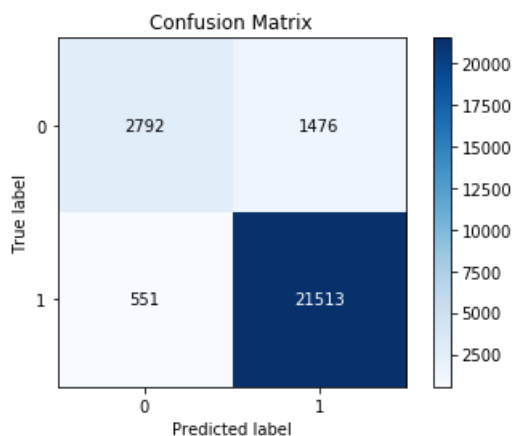
In [56]:

```
#plotting the confusion matrix
#Reference:
#https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html

prediction=clf.predict(final_X_test)
skplt.plot_confusion_matrix(y_test ,prediction)
```

Out[56]:

<matplotlib.axes._subplots.AxesSubplot at 0x9ef4c79908>



[5.2.2] Applying Logistic Regression with L2 regularization on TFIDF, SET 2

In [59]:

```
#Calculating for finding Best C
#predic_proba reference:
#https://stackoverflow.com/questions/37089177/probability-prediction-method-of-kneighborsclassifier-returns-only-0-and-1
#https://discuss.analyticsvidhya.com/t/what-is-the-difference-between-predict-and-predict-proba/67376/3
roc_tr=[]
roc_cv=[]
max_auc_score=0
best_c=0
tuned_parameters=[10**-4, 10**-2, 10**0, 10**2, 10**4]
for i in tuned_parameters:
    clf=LogisticRegression(C=i, penalty='l2', class_weight='balanced')
    # fitting the model on train data
    clf.fit(final_X_tr,y_tr)

    # predict the response on the traininig
    pred_tr = clf.predict_proba(final_X_tr)
    pred_tr=(pred_tr)[:,-1]
    roc_tr.append(roc_auc_score(y_tr,pred_tr))

    # predict the response on the crossvalidation
    pred_cv = clf.predict_proba(final_X_cv)
    pred_cv=(pred_cv)[:,-1]
    roc_cv.append(roc_auc_score(y_cv,pred_cv))

    #finding best c using loop
    if roc_auc_score(y_cv,pred_cv)>max_auc_score:
        best_c=i
        max_auc_score=roc_auc_score(y_cv,pred_cv)

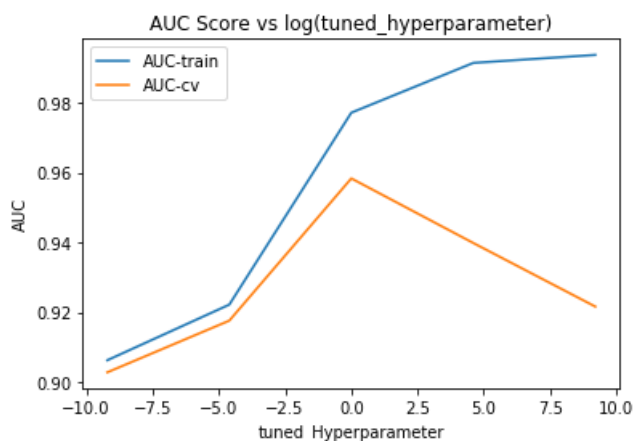
print(best_c)
```

```
print(max_auc_score)
c_l2set2=best_c
auc_set4=max_auc_score
```

```
1
0.9583772216236782
```

In [60]:

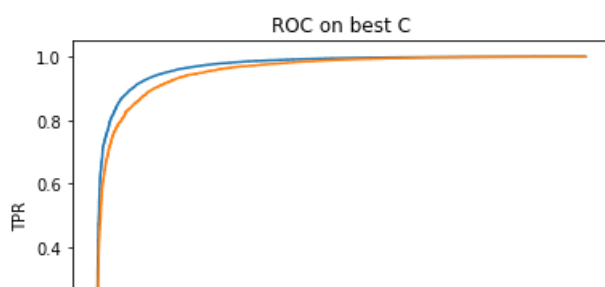
```
# plotting curve between between AUC of cv and train with log of tuned parameter
logalpha=np.log(tuned_parameters)
plt.plot(logalpha,roc_tr,label="AUC-train")
plt.plot(logalpha,roc_cv ,label="AUC-cv")
plt.legend()
plt.xlabel('tuned_Hyperparameter')
plt.ylabel('AUC')
plt.title('AUC Score vs log(tuned_hyperparameter)')
plt.show()
```

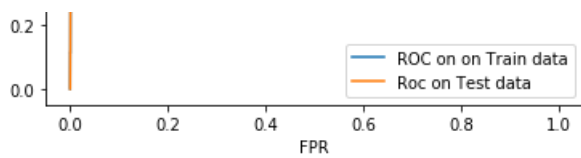


In [61]:

```
#1) Training the model using best C
clf=LogisticRegression(C=c_l2set2, penalty='l2')
clf.fit(final_X_tr, y_tr)
#predicting probability on Test data
pred_test = clf.predict_proba(final_X_test)
pred_test=(pred_test)[:,1]
#predicting probability of Training data
pred_tr = clf.predict_proba(final_X_tr)
pred_tr=(pred_tr)[:,1]

#2)Plotting Roc Curve
#Reference for finding fpr an tpr :
#https://www.programcreek.com/python/example/81207/sklearn.metrics.roc_curve
fpr_tr, tpr_tr, threshold_train = metrics.roc_curve(y_tr, pred_tr)
fpr_test, tpr_test, threshold_test = metrics.roc_curve(y_test, pred_test)
plt.plot(fpr_tr,tpr_tr ,label="ROC on on Train data")
plt.plot(fpr_test,tpr_test ,label="Roc on Test data")
plt.legend()
plt.title('ROC on best C')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```





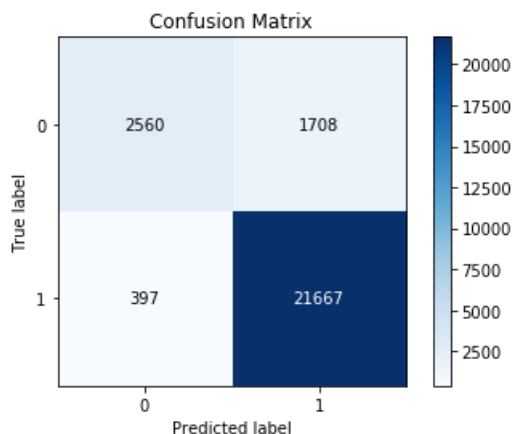
In [62]:

```
#plotting the confusion matrix
#Reference:
#https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html

prediction=clf.predict(final_X_test)
skplt.plot_confusion_matrix(y_test ,prediction)
```

Out[62]:

<matplotlib.axes._subplots.AxesSubplot at 0x9e8760cc18>



[5.2.3] Feature Importance on TFIDF, SET 2

[5.2.3.1] Top 10 important features of positive and negative class from SET 2

In [63]:

```
#Reference for top features is from statistics of machine learning by pratap dangeti:
#https://books.google.co.in/books?id=C-dDDwAAQBAJ&pg=PA216&lpg=PA216&dq=coefs_with_fns%5B:-(n+%2B+
1):-
1&source=bl&ots=j18t1prZXo&sig=ACfU3U2yz8v4v3SOiSrT7fBpJauJKe80DQ&hl=en&sa=X&ved=2ahUKEwjzwK69x6Dj,
XMBHYHqADIQ6AEwAXoECagQAQ#v=onepage&q=coefs_with_fns%5B%3A-(n%20%2B%201)%3A-1&f=false
n=10
feature_names = count_vect.get_feature_names()
coefs_with_fns = sorted(zip(clf.coef_[0], feature_names))
top_n_coefs = zip(coefs_with_fns[:n], coefs_with_fns[:-(n + 1):-1])
print("\tNegative-feature\t\t\t\t\tPositive-feature")
print("-----")
for (coef_1, fn_1), (coef_2, fn_2) in top_n_coefs:
    print("\t%.4f\t%-15s\t\t\t\t\t%.4f\t%-15s" % (coef_1, fn_1, coef_2, fn_2))
```

Negative-feature	Positive-feature		
-6.5446	definite	8.7196	goods
-5.8986	wished	7.1821	beauty
-5.6288	never	7.0413	cube
-5.3333	nursing	5.8382	garden
-5.2658	organic	5.5030	poor
-5.2176	tarts	5.4373	loss
-4.9386	degree	5.2656	local
-4.8892	hole	5.1258	easier
-4.6976	assured	4.9325	whiskey
-4.5900	no	4.7522	always

[5.3] Logistic Regression on AVG W2V, SET 3

[5.3.1] Applying Logistic Regression with L1 regularization on AVG W2V SET 3

In [64]:

```
#Splitting entire data to train,test and cross validation
X=np.array(preprocessed_reviews)
y = np.array(final['Score'])

## split the data set into train and test
X_1, X_test, y_1, y_test = train_test_split(X, y, test_size=0.3, random_state=1)

# split the train data set into cross validation train and cross validation test
X_tr, X_cv, y_tr, y_cv = train_test_split(X_1, y_1, test_size=0.3,random_state=1)

#converting Reviews to Bag of words after splitting to avoid data leakage problem
count_vect = CountVectorizer(min_df=10,max_features=5000)
final_X_tr=count_vect.fit_transform(X_tr)
final_X_test=count_vect.transform(X_test)
final_X_cv=count_vect.transform(X_cv)

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

list_of_sentence_cv=[]
for sentence in X_cv:
    list_of_sentence_cv.append(sentence.split())
final_X_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(list_of_sentence_cv): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to change this
to 300 if you use google's w2v
    cnt_words =0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    final_X_cv.append(sent_vec)

list_of_sentence_test=[]
for sentence in X_test:
    list_of_sentence_test.append(sentence.split())
final_X_test = []; # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(list_of_sentence_test): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to change this
to 300 if you use google's w2v
    cnt_words =0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
```

```
100%|███████████| 43008/43008 [02:57<00:00, 242.85it/s]
100%|███████████| 18433/18433 [01:02<00:00, 295.07it/s]
100%|███████████| 26332/26332 [01:35<00:00, 275.98it/s]
```

```
#Calculating for finding Best C
#predic_proba reference:
#https://stackoverflow.com/questions/37089177/probability-prediction-method-of-kneighborsclassifier-returns-only-0-and-1
#https://discuss.analyticsvidhya.com/t/what-is-the-difference-between-predict-and-predict_proba/67376/3
roc_tr=[]
roc_cv=[]
max_auc_score=0
best_c=0
tuned_parameters=[10**-4, 10**-2, 10**0, 10**2, 10**4]
for i in tuned_parameters:
    clf=LogisticRegression(C=i, penalty='l1', class_weight='balanced')
    # fitting the model on train data
    clf.fit(final_X_tr,y_tr)

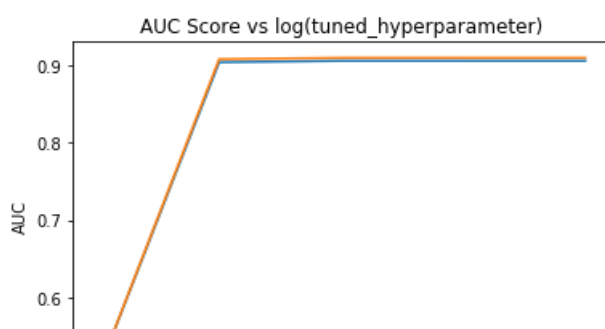
    # predict the response on the traininig
    pred_tr = clf.predict_proba(final_X_tr)
    pred_tr=(pred_tr)[: ,1]
    roc_tr.append(roc_auc_score(y_tr,pred_tr))

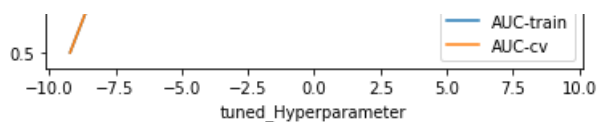
    # predict the response on the crossvalidation
    pred_cv = clf.predict_proba(final_X_cv)
    pred_cv=(pred_cv)[: ,1]
    roc_cv.append(roc_auc_score(y_cv,pred_cv))

#finding best c using loop
if roc_auc_score(y_cv,pred_cv)>max_auc_score:
    best_c=i
    max_auc_score=roc_auc_score(y_cv,pred_cv)

print(best_c)
print(max_auc_score)
c_llset3=best_c
auc set5=max auc score
```

```
# plotting curve between AUC of cv and train with log of tuned parameter
logalpha=np.log(tuned_parameters)
plt.plot(logalpha,roc_tr,label="AUC-train")
plt.plot(logalpha,roc_cv ,label="AUC-cv")
plt.legend()
plt.xlabel('tuned_Hyperparameter')
plt.ylabel('AUC')
plt.title('AUC Score vs log(tuned_hyperparameter)')
plt.show()
```

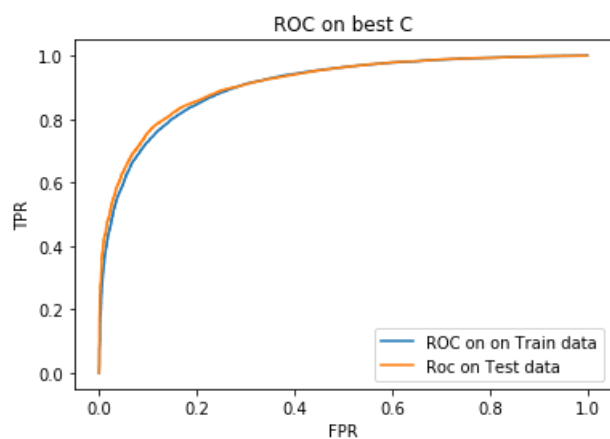




In [67]:

```
#1) Training the model using best C
clf=LogisticRegression(C=c_llset3, penalty='l1')
clf.fit(final_X_tr, y_tr)
#predicting probability on Test data
pred_test = clf.predict_proba(final_X_test)
pred_test=(pred_test)[:,1]
#predicting probability of Training data
pred_tr = clf.predict_proba(final_X_tr)
pred_tr=(pred_tr)[:,1]

#2)Plotting Roc Curve
#Reference for finding fpr an tpr :
#https://www.programcreek.com/python/example/81207/sklearn.metrics.roc_curve
fpr_tr, tpr_tr, threshold_train = metrics.roc_curve(y_tr, pred_tr)
fpr_test, tpr_test, threshold_test = metrics.roc_curve(y_test, pred_test)
plt.plot(fpr_tr,tpr_tr ,label="ROC on on Train data")
plt.plot(fpr_test,tpr_test ,label="Roc on Test data")
plt.legend()
plt.title('ROC on best C')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



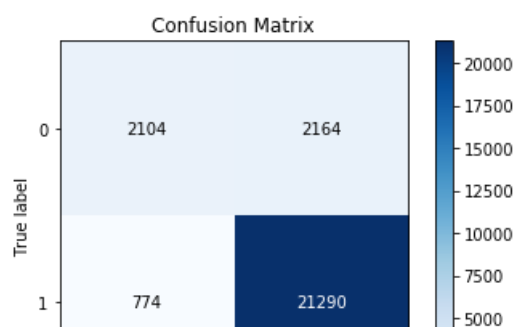
In [68]:

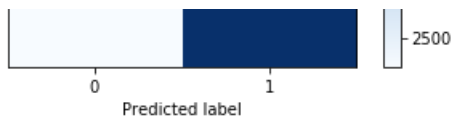
```
#plotting the confusion matrix
#Reference:
#https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html

prediction=clf.predict(final_X_test)
skplt.plot_confusion_matrix(y_test ,prediction)
```

Out[68]:

<matplotlib.axes._subplots.AxesSubplot at 0x9e872d77b8>





[5.3.2] Applying Logistic Regression with L2 regularization on AVG W2V, SET 3

In [69]:

```
#Calculating for finding Best C
#predic_proba reference:
#https://stackoverflow.com/questions/37089177/probability-prediction-method-of-
kneighborsclassifier-returns-only-0-and-1
#https://discuss.analyticsvidhya.com/t/what-is-the-difference-between-predict-and-predict_proba/67
376/3
roc_tr=[]
roc_cv=[]
max_auc_score=0
best_c=0
tuned_parameters=[10**-4, 10**-2, 10**0, 10**2, 10**4]
for i in tuned_parameters:
    clf=LogisticRegression(C=i, penalty='l2', class_weight='balanced')
    # fitting the model on train data
    clf.fit(final_X_tr,y_tr)

    # predict the response on the traininig
    pred_tr = clf.predict_proba(final_X_tr)
    pred_tr=(pred_tr)[: ,1]
    roc_tr.append(roc_auc_score(y_tr,pred_tr))

    # predict the response on the crossvalidation
    pred_cv = clf.predict_proba(final_X_cv)
    pred_cv=(pred_cv)[: ,1]
    roc_cv.append(roc_auc_score(y_cv,pred_cv))

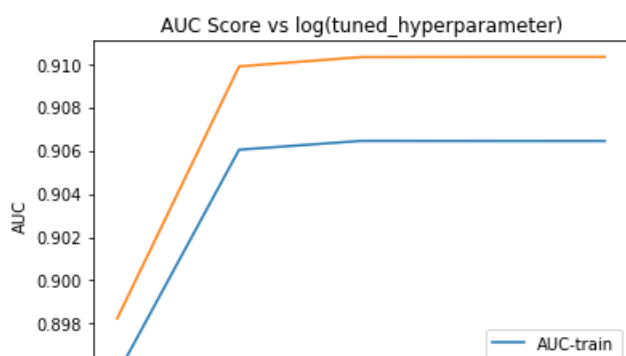
    #finding best c using loop
    if roc_auc_score(y_cv,pred_cv)>max_auc_score:
        best_c=i
        max_auc_score=roc_auc_score(y_cv,pred_cv)

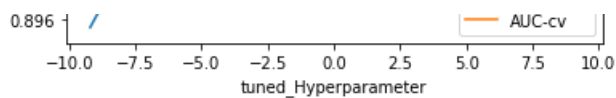
print(best_c)
print(max_auc_score)
c_l2set3=best_c
auc_set6=max_auc_score
```

10000
0.9103343224622573

In [70]:

```
# plotting curve between between AUC of cv and train with log of tuned parameter
plt.plot(np.log(tuned_parameters),roc_tr,label="AUC-train")
plt.plot(np.log(tuned_parameters),roc_cv ,label="AUC-cv")
plt.legend()
plt.xlabel('tuned_Hyperparameter')
plt.ylabel('AUC')
plt.title('AUC Score vs log(tuned_hyperparameter) ')
plt.show()
```

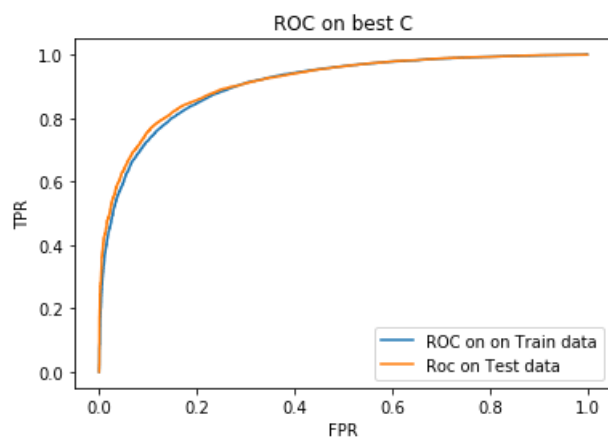




In [71]:

```
#1) Training the model using best C
clf=LogisticRegression(C=c_l2set3, penalty='l2')
clf.fit(final_X_tr, y_tr)
#predicting probability on Test data
pred_test = clf.predict_proba(final_X_test)
pred_test=(pred_test)[:,1]
#predicting probability of Training data
pred_tr = clf.predict_proba(final_X_tr)
pred_tr=(pred_tr)[:,1]

#2)Plotting Roc Curve
#Reference for finding fpr an tpr :
#https://www.programcreek.com/python/example/81207/sklearn.metrics.roc_curve
fpr_tr, tpr_tr, threshold_train = metrics.roc_curve(y_tr, pred_tr)
fpr_test, tpr_test, threshold_test = metrics.roc_curve(y_test, pred_test)
plt.plot(fpr_tr,tpr_tr ,label="ROC on on Train data")
plt.plot(fpr_test,tpr_test ,label="Roc on Test data")
plt.legend()
plt.title('ROC on best C')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



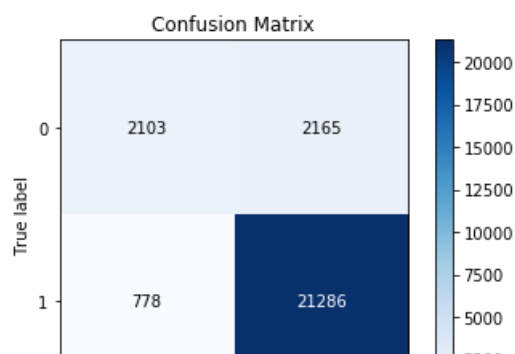
In [72]:

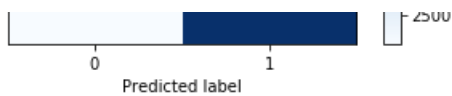
```
#plotting the confusion matrix
#Reference:
#https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html

prediction=clf.predict(final_X_test)
skplt.plot_confusion_matrix(y_test ,prediction)
```

Out[72]:

<matplotlib.axes._subplots.AxesSubplot at 0x9e933c3d30>





[5.4] Logistic Regression on TFIDF W2V, SET 4

[5.4.1] Applying Logistic Regression with L1 regularization on TFIDF W2V, SET 4

In [73]:

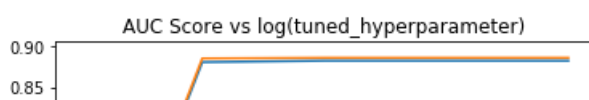
```
# TF-IDF weighted Word2Vec
tfidf_feat = model.get_feature_names() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val = tfidf

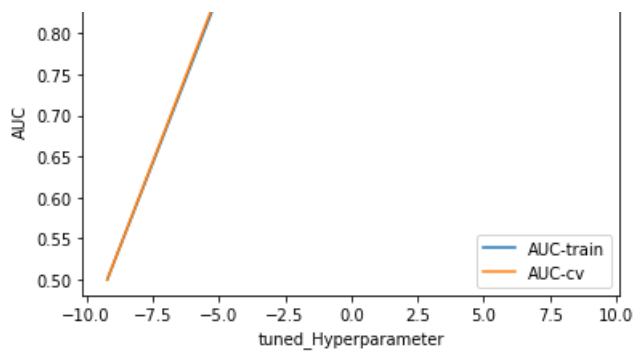
list_of_sentence_tr=[]
for sentence in X_tr:
    list_of_sentence_tr.append(sentence.split())
final_X_tr = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(list_of_sentence_tr): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of 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 courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    final_X_tr.append(sent_vec)
    row += 1

list_of_sentence_cv=[]
for sentence in X_cv:
    list_of_sentence_cv.append(sentence.split())
final_X_cv = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(list_of_sentence_cv): # 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 courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    final_X_cv.append(sent_vec)
    row += 1

list_of_sentence_test=[]
for sentence in X_test:
    list_of_sentence_test.append(sentence.split())
final_X_test = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(list_of_sentence_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]
```

```
100%|██████████| 43008/43008 [04:44<00:00, 151.20it/s]
100%|██████████| 18433/18433 [02:06<00:00, 146.26it/s]
100%|██████████| 26332/26332 [03:06<00:00, 140.93it/s]
```

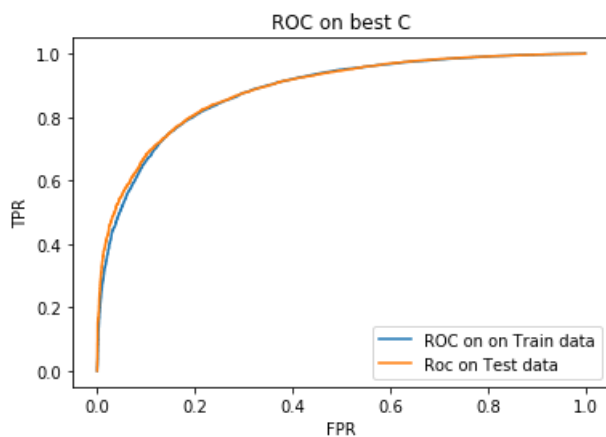




In [77]:

```
#1) Training the model using best C
clf=LogisticRegression(C=c_llset4, penalty='l1')
clf.fit(final_X_tr, y_tr)
#predicting probability on Test data
pred_test = clf.predict_proba(final_X_test)
pred_test=(pred_test)[:,1]
#predicting probability of Training data
pred_tr = clf.predict_proba(final_X_tr)
pred_tr=(pred_tr)[:,1]

#2)Plotting Roc Curve
#Reference for finding fpr an tpr :
#https://www.programcreek.com/python/example/81207/sklearn.metrics.roc_curve
fpr_tr, tpr_tr, threshold_train = metrics.roc_curve(y_tr, pred_tr)
fpr_test, tpr_test, threshold_test = metrics.roc_curve(y_test, pred_test)
plt.plot(fpr_tr,tpr_tr ,label="ROC on on Train data")
plt.plot(fpr_test,tpr_test ,label="Roc on Test data")
plt.legend()
plt.title('ROC on best C')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



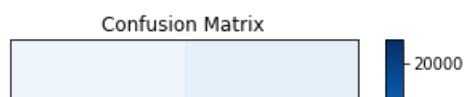
In [78]:

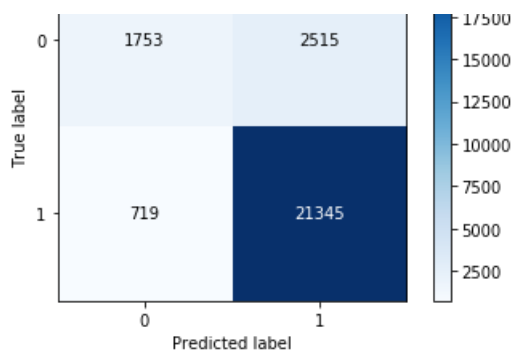
```
#plotting the confusion matrix
#Reference:
#https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html

prediction=clf.predict(final_X_test)
skplt.plot_confusion_matrix(y_test ,prediction)
```

Out[78]:

<matplotlib.axes._subplots.AxesSubplot at 0x9ea7c36978>





[5.4.2] Applying Logistic Regression with L2 regularization on TFIDF W2V, SET 4

In [79]:

```
#Calculating for finding Best C
#predic_proba reference:
#https://stackoverflow.com/questions/37089177/probability-prediction-method-of-
#kneighborsclassifier-returns-only-0-and-1
#https://discuss.analyticsvidhya.com/t/what-is-the-difference-between-predict-and-predict-proba/67
#376/3
roc_tr=[]
roc_cv=[]
max_auc_score=0
best_c=0
tuned_parameters=[10**-4, 10**-2, 10**0, 10**2, 10**4]
for i in tuned_parameters:
    clf=LogisticRegression(C=i, penalty='l2', class_weight='balanced')
    # fitting the model on train data
    clf.fit(final_X_tr,y_tr)

    # predict the response on the traininig
    pred_tr = clf.predict_proba(final_X_tr)
    pred_tr=(pred_tr)[: ,1]
    roc_tr.append(roc_auc_score(y_tr,pred_tr))

    # predict the response on the crossvalidation
    pred_cv = clf.predict_proba(final_X_cv)
    pred_cv=(pred_cv)[: ,1]
    roc_cv.append(roc_auc_score(y_cv,pred_cv))

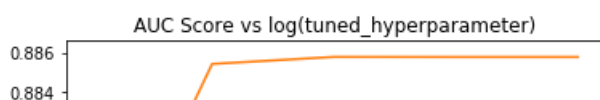
    #finding best c using loop
    if roc_auc_score(y_cv,pred_cv)>max_auc_score:
        best_c=i
        max_auc_score=roc_auc_score(y_cv,pred_cv)

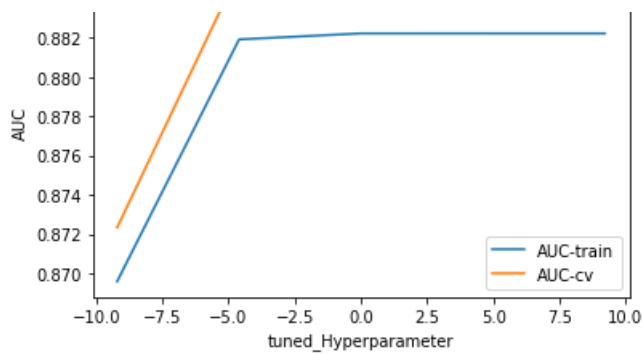
print(best_c)
print(max_auc_score)
c_l2set4=best_c
auc_set8=max_auc_score
```

1
0.8857726189154723

In [80]:

```
# plotting curve between between AUC of cv and train with log of tuned parameter
logalpha=np.log(tuned_parameters)
plt.plot(logalpha,roc_tr,label="AUC-train")
plt.plot(logalpha,roc_cv ,label="AUC-cv")
plt.legend()
plt.xlabel('tuned_Hyperparameter')
plt.ylabel('AUC')
plt.title('AUC Score vs log(tuned_hyperparameter)')
plt.show()
```

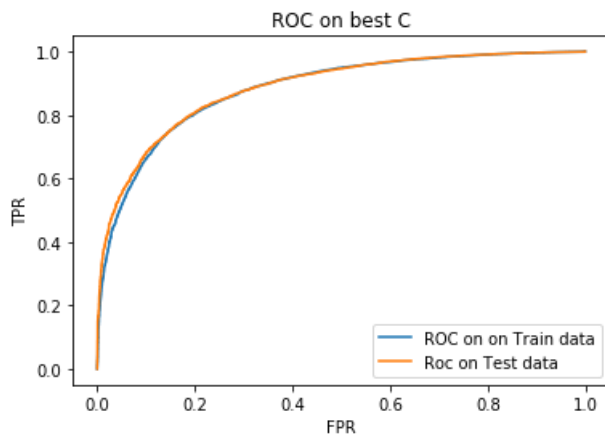




In [81]:

```
#1) Training the model using best C
clf=LogisticRegression(C=c_l2set4, penalty='l1')
clf.fit(final_X_tr, y_tr)
#predicting probability on Test data
pred_test = clf.predict_proba(final_X_test)
pred_test=(pred_test)[:,1]
#predicting probability of Training data
pred_tr = clf.predict_proba(final_X_tr)
pred_tr=(pred_tr)[:,1]

#2)Plotting Roc Curve
#Reference for finding fpr an tpr :
#https://www.programcreek.com/python/example/81207/sklearn.metrics.roc_curve
fpr_tr, tpr_tr, threshold_train = metrics.roc_curve(y_tr, pred_tr)
fpr_test, tpr_test, threshold_test = metrics.roc_curve(y_test, pred_test)
plt.plot(fpr_tr,tpr_tr ,label="ROC on on Train data")
plt.plot(fpr_test,tpr_test ,label="Roc on Test data")
plt.legend()
plt.title('ROC on best C')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



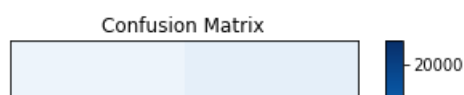
In [82]:

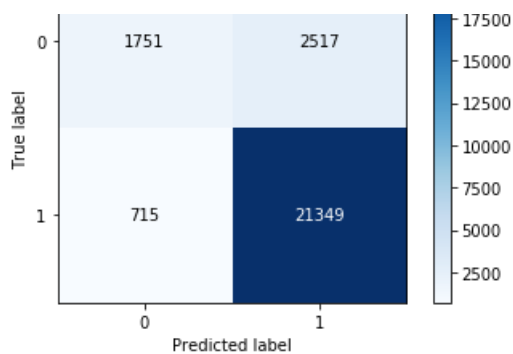
```
#plotting the confusion matrix
#Reference:
#https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html

prediction=clf.predict(final_X_test)
skplt.plot_confusion_matrix(y_test ,prediction)
```

Out[82]:

<matplotlib.axes._subplots.AxesSubplot at 0x9ea7c20048>





[6] Conclusions

In [83]:

```
# Please compare all your models using Prettytable library
from prettytable import PrettyTable
x = PrettyTable()
x.field_names = ["Vectorizer", "Model_name", "Hyperamete-best_k", "AUC"]
x.add_row(["BOW", "L1", c_llset1, auc_set1])
x.add_row(["TFIDF", "L2", c_l2set1, auc_set2])
x.add_row(["AwgW2V", "L1", c_llset2, auc_set3])
x.add_row(["TFIDF-W2V", "L2", c_l2set2, auc_set4])
x.add_row(["BOW", "L1", c_llset3, auc_set5])
x.add_row(["TFIDF", "L2", c_l2set3, auc_set6])
x.add_row(["AwgW2V", "L1", c_llset4, auc_set7])
x.add_row(["TFIDF-W2V", "L2", c_l2set4, auc_set8])
print(x)
```

Vectorizer	Model_name	Hyperamete-best_k	AUC
BOW	L1	1	0.9247882312415648
TFIDF	L2	0.01	0.9318337422621068
AwgW2V	L1	1	0.9558029952014784
TFIDF-W2V	L2	1	0.9583772216236782
BOW	L1	1	0.9103362393503098
TFIDF	L2	10000	0.9103343224622573
AwgW2V	L1	10000	0.8857688732721511
TFIDF-W2V	L2	1	0.8857726189154723