

Amazon Fine Food Reviews Analysis

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

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

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

Number of reviews: 568,454

Number of users: 256,059

Number of products: 74,258

Timespan: Oct 1999 - Oct 2012

Number of Attributes/Columns in data: 10

Attribute Information:

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

Objective:

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

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

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

[1]. Reading Data

[1.1] Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

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

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

In [1]:

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

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

```

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

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

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

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

from tqdm import tqdm
import os

```

In [2]:

```

# using SQLite Table to read data.
con = sqlite3.connect('C:/Downloads/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 5000""", 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 (5000, 10)

Out[2]:

	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
2	3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia	1	1	1	1219017600	"Delight" says it all

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary
--	----	-----------	--------	-------------	----------------------	------------------------	-------	------	---------

In [3]:

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

In [4]:

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

(80668, 7)

Out[4]:

	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 [5]:

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

Out[5]:

	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 [6]:

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

Out[6]:

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 [7]:

```
display= pd.read_sql_query("""
SELECT *
FROM Reviews
```

```
FROM reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID
""", con)
display.head()
```

Out[7]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summ
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA ⁺ VANII WAFE
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA ⁺ VANII WAFE
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA ⁺ VANII WAFE
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA ⁺ VANII WAFE
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA ⁺ 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 [8]:

```
#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 [9]:

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

Out[9]:

(4986, 10)

In [10]:

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

Out[10]:

99.72

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 [11]:

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

display.head()
```

Out[11]:

	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 [12]:

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

In [13]:

```
#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()
```

(4986, 10)

Out[13]:

```
1    4178
0     808
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 [14]:

```
# 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)
```

Why is this \$[...] when the same product is available for \$[...] here?
<http://www.amazon.com/VICTOR-FLY-MAGNET-BAIT-REFILL/dp/B00004RBDY>

The Victor M380 and M502 traps are unreal, of course -- total fly genocide. Pretty stinky, but only right nearby.

=====

I recently tried this flavor/brand and was surprised at how delicious these chips are. The best thing was that there were a lot of "brown" chips in the bsg (my favorite), so I bought some more through amazon and shared with family and friends. I am a little disappointed that there are not, so far, very many brown chips in these bags, but the flavor is still very good. I like them better than the yogurt and green onion flavor because they do not seem to be as salty, and the onion flavor is better. If you haven't eaten Kettle chips before, I recommend that you try a bag before buying bulk. They are thicker and crunchier than Lays but just as fresh out of the bag.

=====

Wow. So far, two two-star reviews. One obviously had no idea what they were ordering; the other wants crispy cookies. Hey, I'm sorry; but these reviews do nobody any good beyond reminding us to look before ordering.

These are chocolate-oatmeal cookies. If you don't like that combination, don't order this type of cookie. I find the combo quite nice, really. The oatmeal sort of "calms" the rich chocolate flavor and gives the cookie sort of a coconut-type consistency. Now let's also remember that tastes differ; so, I've given my opinion.

Then, these are soft, chewy cookies -- as advertised. They are not "crispy" cookies, or the blurb would say "crispy," rather than "chewy." I happen to like raw cookie dough; however, I don't see where these taste like raw cookie dough. Both are soft, however, so is this the confusion? And, yes, they stick together. Soft cookies tend to do that. They aren't individually wrapped, which would add to the cost. Oh yeah, chocolate chip cookies tend to be somewhat sweet.

So, if you want something hard and crisp, I suggest Nabisco's Ginger Snaps. If you want a cookie that's soft, chewy and tastes like a combination of chocolate and oatmeal, give these a try. I'm here to place my second order.

=====

love to order my coffee on amazon. easy and shows up quickly.
This k cup is great coffee. dcafe is very good as well

In [15]:

```
# 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_1500 = re.sub(r"http\S+", "", sent_1500)
sent_4900 = re.sub(r"http\S+", "", sent_4900)

print(sent_0)
```

Why is this \$[...] when the same product is available for \$[...] here?

The Victor M380 and M502 traps are unreal, of course -- total fly genocide. Pretty stinky, but only right nearby.

In [16]:

In [10]:

```
# 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)
```

Why is this \$[...] when the same product is available for \$[...] here? />The Victor M380 and M502 traps are unreal, of course -- total fly genocide. Pretty stinky, but only right nearby.

=====

I recently tried this flavor/brand and was surprised at how delicious these chips are. The best thing was that there were a lot of "brown" chips in the bsg (my favorite), so I bought some more through amazon and shared with family and friends. I am a little disappointed that there are not, so far, very many brown chips in these bags, but the flavor is still very good. I like them better than the yogurt and green onion flavor because they do not seem to be as salty, and the onion flavor is better. If you haven't eaten Kettle chips before, I recommend that you try a bag before buying bulk. They are thicker and crunchier than Lays but just as fresh out of the bag.

=====

Wow. So far, two two-star reviews. One obviously had no idea what they were ordering; the other wants crispy cookies. Hey, I'm sorry; but these reviews do nobody any good beyond reminding us to look before ordering. These are chocolate-oatmeal cookies. If you don't like that combination, don't order this type of cookie. I find the combo quite nice, really. The oatmeal sort of "calms" the rich chocolate flavor and gives the cookie sort of a coconut-type consistency. Now let's also remember that tastes differ; so, I've given my opinion. Then, these are soft, chewy cookies -- as advertised. They are not "crispy" cookies, or the blurb would say "crispy," rather than "chewy." I happen to like raw cookie dough; however, I don't see where these taste like raw cookie dough. Both are soft, however, so is this the confusion? And, yes, they stick together. Soft cookies tend to do that. They aren't individually wrapped, which would add to the cost. Oh yeah, chocolate chip cookies tend to be somewhat sweet. So, if you want something hard and crisp, I suggest Nabisco's Ginger Snaps. If you want a cookie that's soft, chewy and tastes like a combination of chocolate and oatmeal, give these a try. I'm here to place my second order.

=====

love to order my coffee on amazon. easy and shows up quickly. This 1/2 cup is great coffee. decaf is very good as well

In [17]:

```
# 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 [18]:

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

Wow. So far, two two-star reviews. One obviously had no idea what they were ordering; the other wants crispy cookies. Hey, I am sorry; but these reviews do nobody any good beyond reminding us to look before ordering.

These are chocolate-oatmeal cookies. If you do not like that combination, do not order this type of cookie. I find the combo quite nice, really. The oatmeal sort of "calms" the rich chocolate flavor and gives the cookie sort of a coconut-type consistency. Now let is also remember that tastes differ; so, I have given my opinion.

Then, these are soft, chewy cookies -- as advertised. They are not "crispy" cookies, or the blurb would say "crispy," rather than "chewy." I happen to like raw cookie dough; however, I do not see where these taste like raw cookie dough. Both are soft, however, so is this the confusion? And, yes, they stick together. Soft cookies tend to do that. They are not individually wrapped, which would add to the cost. Oh yeah, chocolate chip cookies tend to be somewhat sweet.

So, if you want something hard and crisp, I suggest Nabisco is Ginger Snaps. If you want a cookie that is soft, chewy and tastes like a combination of chocolate and oatmeal, give these a try. I am here to place my second order.

=====

In [19]:

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

Why is this \$[...] when the same product is available for \$[...] here?
 />
The Victor and traps are unreal, of course -- total fly genocide. Pretty stinky, but only right nearby.

In [20]:

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

Wow So far two two star reviews One obviously had no idea what they were ordering the other wants crispy cookies Hey I am sorry but these reviews do nobody any good beyond reminding us to look before ordering br br These are chocolate oatmeal cookies If you do not like that combination do not order this type of cookie I find the combo quite nice really The oatmeal sort of calms the rich chocolate flavor and gives the cookie sort of a coconut type consistency Now let is also remember that tastes differ so I have given my opinion br br Then these are soft chewy cookies as advertised They are not crispy cookies or the blurb would say crispy rather than chewy I happen to like raw cookie dough however I do not see where these taste like raw cookie dough Both are soft however so is this the confusion And yes they stick together Soft cookies tend to do that They are not individually wrapped which would add to the cost Oh yeah chocolate chip cookies tend to be somewhat sweet br br So if you want something hard and crisp I suggest Nabisco is Ginger Snaps If you want a cookie that is soft chewy and tastes like a combination of chocolate and oatmeal give these a try I am here to place my second order

In [21]:

```
# 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', "you'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', \
```


In [22]:

```
100%|███████████| 4986/4986 [00:03<00:00, 1341.75it/s]
```

In [23]:

Out [23] :

[3.2] Preprocessing Review Summary

In [24]:

[4] Featurization

[4.1] BAG OF WORDS

In [25]:

```
#BoW
count_vect = CountVectorizer() #in scikit-learn
count_vect.fit(preprocessed_reviews)
print("some feature names ", count_vect.get_feature_names()[: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  ['aa', 'aahhhs', 'aback', 'abandon', 'abates', 'abbott', 'abby', 'abdominal',
'abiding', 'ability']
=====
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (4986, 12997)
the number of unique words  12997

```

[4.2] Bi-Grams and n-Grams.

In [26]:

```

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

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

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

```

```

the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (4986, 3144)
the number of unique words including both unigrams and bigrams  3144

```

[4.3] TF-IDF

In [27]:

```

tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
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)

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 find', 'able get',
'absolute', 'absolutely', 'absolutely delicious', 'absolutely love', 'absolutely no', 'according']
=====
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer  (4986, 3144)
the number of unique words including both unigrams and bigrams  3144

```

[4.4] Word2Vec

In [28]:

```

# 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 [29]:

```
# 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

is_your_ram_gt_16g=False
want_to_use_google_w2v = False
want_to_train_w2v = True

if want_to_train_w2v:
    # min_count = 5 considers only words that 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'))

elif want_to_use_google_w2v and is_your_ram_gt_16g:
    if os.path.isfile('GoogleNews-vectors-negative300.bin'):
        w2v_model=KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin', binary=True)
        print(w2v_model.wv.most_similar('great'))
        print(w2v_model.wv.most_similar('worst'))
    else:
        print("you don't have google's word2vec file, keep want_to_train_w2v = True, to train your own w2v ")

[('greasy', 0.9928233623504639), ('salty', 0.9927825927734375), ('tasty', 0.9927698969841003),
 ('flavorful', 0.9927244186401367), ('subtle', 0.992603063583374), ('delicious',
 0.9922800660133362), ('alternative', 0.9922574758529663), ('texture', 0.9920745491981506),
 ('want', 0.9919319748878479), ('enjoy', 0.9917845726013184)]
=====
[('style', 0.9992372989654541), ('choice', 0.9992095232009888), ('surprised', 0.9992057085037231),
 ('similar', 0.9991554021835327), ('stand', 0.9991458058357239), ('drop', 0.9991284012794495),
 ('usual', 0.9991222023963928), ('type', 0.9990949034690857), ('remember', 0.9990780353546143), ('p
opcorn', 0.9990752339363098)]
```

In [30]:

```
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 3817
sample words ['product', 'available', 'course', 'total', 'pretty', 'stinky', 'right', 'nearby', '
used', 'ca', 'not', 'beat', 'great', 'received', 'shipment', 'could', 'hardly', 'wait', 'try', 'lo
ve', 'call', 'instead', 'removed', 'easily', 'daughter', 'designed', 'printed', 'use', 'car', 'win
dows', 'beautifully', 'shop', 'program', 'going', 'lot', 'fun', 'everywhere', 'like', 'tv',
'computer', 'really', 'good', 'idea', 'final', 'outstanding', 'window', 'everybody', 'asks',
'bought', 'made']
```

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

[4.4.1.1] Avg W2v

In [31]:

```
# average Word2Vec
```

[illegible][illegible]

- **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. Apply Knn(kd tree version) on these feature sets

NOTE: sklearn implementation of kd-tree accepts only dense matrices, you need to convert the sparse matrices of CountVectorizer/TfidfVectorizer into dense matrices. You can convert sparse matrices to dense using `.toarray()` attribute. For more information please visit this [link](#)

- **SET 5:** Review text, preprocessed one converted into vectors using (BOW) but with restriction on maximum features generated.

```
count_vect = CountVectorizer(min_df=10, max_features=500)
count_vect.fit(preprocessed_reviews)
```

- **SET 6:** Review text, preprocessed one converted into vectors using (TFIDF) but with restriction on maximum features generated.

```
tf_idf_vect = TfidfVectorizer(min_df=10, max_features=500)
tf_idf_vect.fit(preprocessed_reviews)
```

- **SET 3:** Review text, preprocessed one converted into vectors using (AVG W2v)
- **SET 4:** Review text, preprocessed one converted into vectors using (TFIDF W2v)

3. The hyper paramter tuning(find best K)

- Find the best hyper parameter which will give the maximum [AUC](#) value
- Find the best hyper paramter using k-fold cross validation or simple cross validation data
- Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

4. 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](#) with predicted and original labels of test data points

5. Conclusion

- You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this prettytable library [link](#)

Note: Data Leakage

1. There will be an issue of data-leakage if you vectorize the entire data and then split it into train/cv/test.
2. To avoid the issue of data-leakage, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method `fit_transform()` on you train data, and apply the method `transform()` on cv/test data.
4. For more details please go through this [link](#).

[5.1] Applying KNN brute force

[5.1.1] Applying KNN brute force on BOW, **SET 1**

5.1.11 Importing required libraries

In [34]:

```
# ===== loading libraries =====
import pdb
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import sklearn
```

```

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
# =====

```

5.1.12 Splitting the data converting to bag of words

In [35]:

```

#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()
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)

```

5.1.13 Hyper parameter tuning-Finding the best k using simple cross validation

In [36]:

```

#Calculating for finding Best K
#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
k_value=[]
roc_tr=[]
roc_cv=[]
max_auc_score=0
best_k=0
for i in range(1,50,2):
    # instantiate learning model (k = 50)
    knn = KNeighborsClassifier(n_neighbors=i,algorithm='brute',metric='minkowski')

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

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

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

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

print(best_k)
print(max_auc_score)
k_set1=best_k
auc_set1=max_auc_score

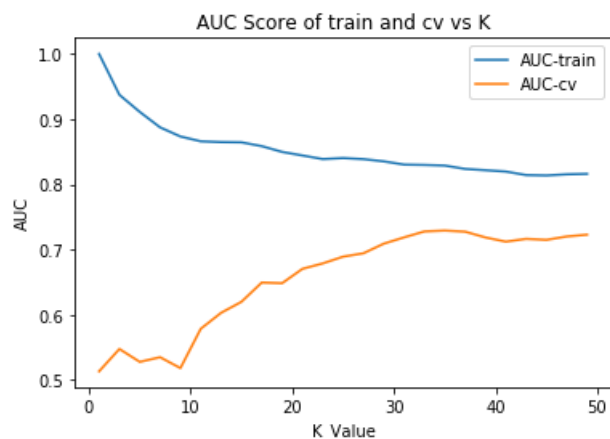
```

35
0.7291652061747265

5.1.14 Curve plotting between AUC of cv and train with k

In [37]:

```
# plotting curve between between AUC of cv and train with k
plt.plot(k_value,roc_tr,label="AUC-train")
plt.plot(k_value,roc_cv ,label="AUC-cv")
plt.legend()
plt.xlabel('K_Value')
plt.ylabel('AUC')
plt.title('AUC Score of train and cv vs K')
plt.show()
```

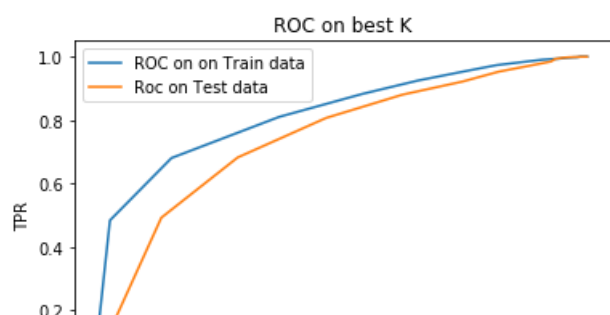


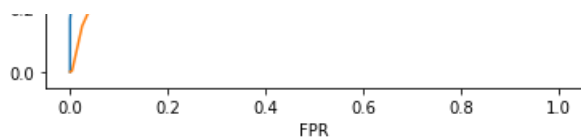
5.1.15 Training the model with the obtained best_k and plotting Roc curve

In [38]:

```
#1) Training the model using best K
knn = KNeighborsClassifier(n_neighbors=best_k,algorithm='brute',metric='minkowski')
knn.fit(final_X_tr, y_tr)
#predicting probability of success on Training data
pred_tr = knn.predict_proba(final_X_tr)
pred_tr=(pred_tr)[:,1]
#predicting probability of success on Test data
pred_test = knn.predict_proba(final_X_test)
pred_test=(pred_test)[:,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 K')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```





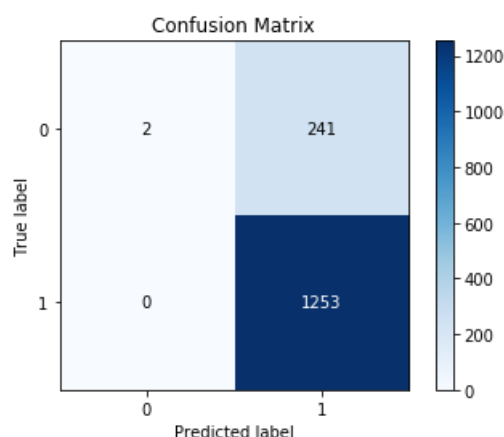
In [39]:

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

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

Out[39]:

<matplotlib.axes._subplots.AxesSubplot at 0x10d582f0>



[5.1.2] Applying KNN brute force on TFIDF, SET 2

5.1.21 Splitting the data converting to TFIDF

In [40]:

```
#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
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2),min_df=10 )
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)
```

5.1.22 Hyper parameter tuning-Finding the best k using simple cross validation

In [41]:

```
#Calculating for finding Best K
#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
k_value=[]
roc_tr=[]
```



```

roc_cv=[]
max_auc_score=0
best_k=0
for i in range(1,50,2):
    # instantiate learning model (k = 50)
    knn = KNeighborsClassifier(n_neighbors=i,algorithm='brute',metric='minkowski')

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

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

    # predict the response on the training
    pred_tr = knn.predict_proba(final_X_tr)
    pred_tr=(pred_tr)[:,1]
    roc_tr.append(roc_auc_score(y_tr,pred_tr))
    k_value.append(i)

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

print(best_k)
print(max_auc_score)
k_set2=best_k
auc_set2=max_auc_score

```

45
0.8345484626277638

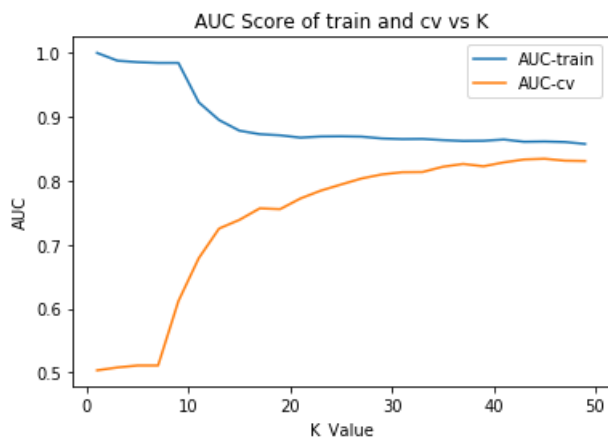
5.1.23 Curve plotting between AUC of cv and train with k

In [42]:

```

# plotting curve between between AUC of cv and train with k
plt.plot(k_value,roc_tr,label="AUC-train")
plt.plot(k_value,roc_cv ,label="AUC-cv")
plt.legend()
plt.xlabel('K_Value')
plt.ylabel('AUC')
plt.title('AUC Score of train and cv vs K')
plt.show()

```



5.1.24 Training the model with the obtained best_k and plotting Roc curve

In [43]:

```

#1) Training the model using best K
knn = KNeighborsClassifier(n_neighbors=best_k,algorithm='brute',metric='minkowski')
knn.fit(final_X_tr, y_tr)
#predicting probability of success Training data

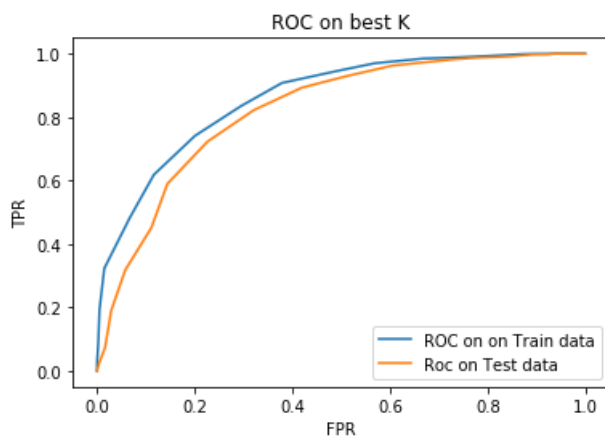
```

```

pred_tr = knn.predict_proba(final_X_tr)
pred_tr=(pred_tr)[: ,1]
#predicting probability of success on Test data
pred_test = knn.predict_proba(final_X_test)
pred_test=(pred_test)[: ,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 K')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()

```



In [44]:

```

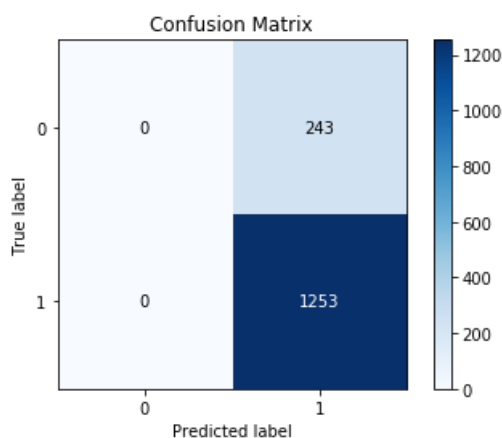
#plotting the confusion matrix
#Reference:
#https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html

prediction=knn.predict(final_X_test)
skplt.plot_confusion_matrix(y_test ,prediction)

```

Out[44]:

<matplotlib.axes._subplots.AxesSubplot at 0x106c3bb0>



[5.1.3] Applying KNN brute force on AVG W2V, SET 3

In [45]:

```

#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)

# Please write all the code with proper documentation
# 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
    if cnt_words != 0:
        sent_vec /= cnt_words
    final_X_test.append(sent_vec)

```

```

100%|████████████████████████████████████████| 2443/2443 [00:05<00:00, 447.74it/s]
100%|████████████████████████████████████████| 1047/1047 [00:02<00:00, 452.63it/s]
100%|████████████████████████████████████████| 1496/1496 [00:03<00:00, 445.21it/s]

```

In [97]:

```

#Calculating for finding Best K
#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
k_value=[]
roc_tr=[]
roc_cv=[]
max_auc_score=0
best_k=0
for i in range(1,50,2):
    # instantiate learning model (k = 50)
    knn = KNeighborsClassifier(n_neighbors=i,algorithm='brute',metric='minkowski')

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

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

    # predict the response on the training
    pred_tr = knn.predict_proba(final_X_tr)
    pred_tr=(pred_tr)[:,1]
    roc_tr.append(roc_auc_score(y_tr,pred_tr))
    k_value.append(i)

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

print(best_k)
print(max_auc_score)
k_set3=best_k
auc_set3=max_auc_score

```

```

41
0.6050993368198198

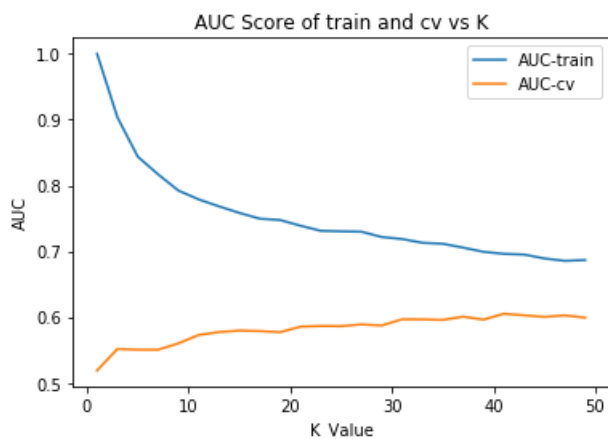
```

In [47]:

```

# plotting curve between between AUC of cv and train with k
plt.plot(k_value,roc_tr,label="AUC-train")
plt.plot(k_value,roc_cv ,label="AUC-cv")
plt.legend()
plt.xlabel('K_Value')
plt.ylabel('AUC')
plt.title('AUC Score of train and cv vs K')
plt.show()

```



In [48]:

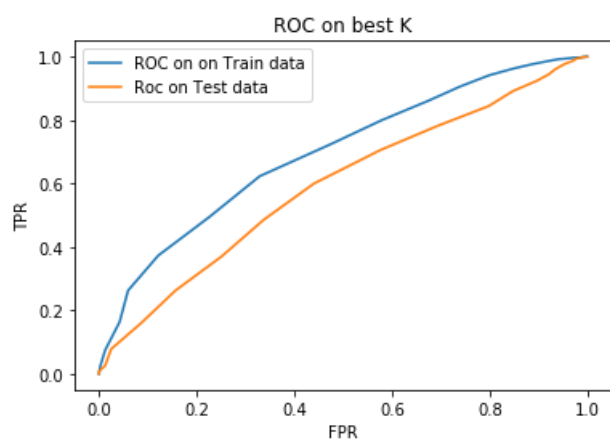
```

#1) Training the model using best K
knn = KNeighborsClassifier(n_neighbors=best_k,algorithm='brute',metric='minkowski')
knn.fit(final_X_tr, y_tr)
#predicting probability of success Training data
pred_tr = knn.predict_proba(final_X_tr)
pred_tr=(pred_tr)[:,1]
#predicting probability of success Test data

```

```
#predicting probability of success on test data
pred_test = knn.predict_proba(final_X_test)
pred_test=(pred_test)[: ,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 K')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



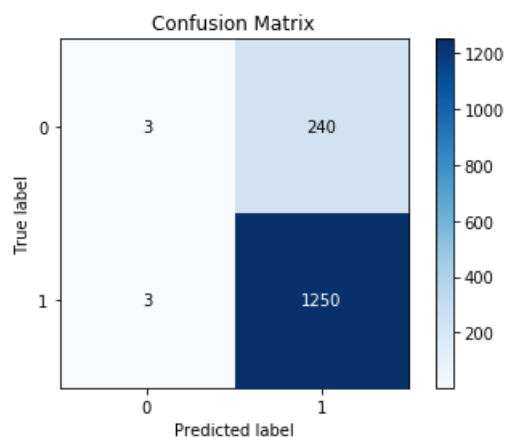
In [49]:

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

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

Out[49]:

<matplotlib.axes._subplots.AxesSubplot at 0x1174c790>



[5.1.4] Applying KNN brute force on TFIDF W2V, SET 4

In [50]:

```
# 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 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
    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 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
    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]
            #
            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
    final_X_test.append(sent_vec)
    row += 1

```

```

100%|████████████████████████████████████████| 2443/2443 [00:26<00:00, 92.51it/s]
100%|████████████████████████████████████████| 1047/1047 [00:12<00:00, 86.97it/s]
100%|████████████████████████████████████████| 1496/1496 [00:16<00:00, 92.51it/s]

```

In [98]:

```
#Calculating for finding Best K
#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
k_value=[]
roc_tr=[]
roc_cv=[]
max_auc_score=0
best_k=0
for i in range(1,50,2):
    # instantiate learning model (k = 50)
    knn = KNeighborsClassifier(n_neighbors=i,algorithm='brute',metric='minkowski')

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

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

    # predict the response on the training
    pred_tr = knn.predict_proba(final_X_tr)
    pred_tr=(pred_tr)[:,1]
    roc_tr.append(roc_auc_score(y_tr,pred_tr))
    k_value.append(i)

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

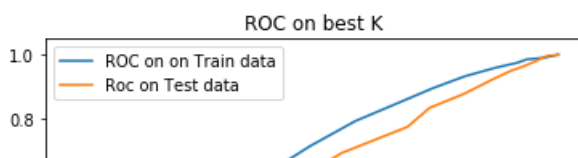
print(best_k)
print(max_auc_score)
k_set4=best_k
auc_set4=max_auc_score
```

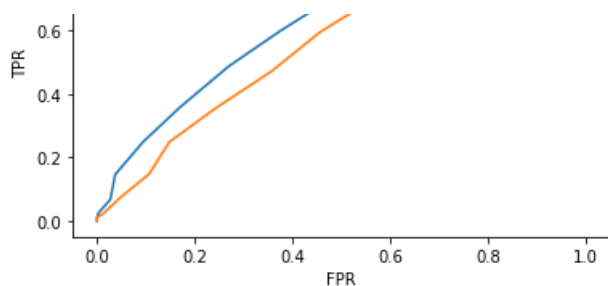
41
0.6050993368198198

In [52]:

```
#1)Training the model using best K
knn = KNeighborsClassifier(n_neighbors=best_k,algorithm='brute',metric='minkowski')
knn.fit(final_X_tr, y_tr)
#predicting probability of success Training data
pred_tr = knn.predict_proba(final_X_tr)
pred_tr=(pred_tr)[:,1]
#predicting probability of success on Test data
pred_test = knn.predict_proba(final_X_test)
pred_test=(pred_test)[:,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 K')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```





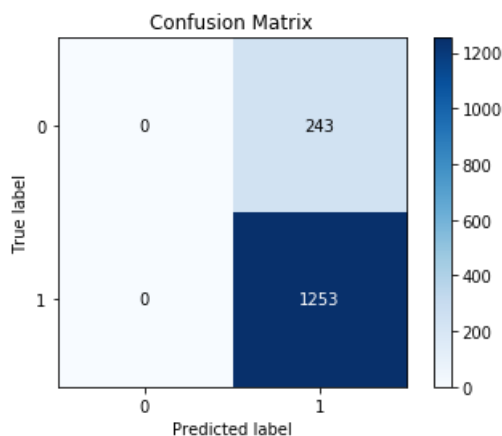
In [53]:

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

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

Out[53]:

<matplotlib.axes._subplots.AxesSubplot at 0x1180c490>



[5.2] Applying KNN kd-tree

[5.2.1] Applying KNN kd-tree on BOW, SET 5

In [54]:

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

#To over come memmory error issue more splitting has been done
X, X_memoryerror, y, y_memoryerror = train_test_split(X, y, test_size=0.70, random_state=1)

# 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()
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)

#converting them to dense because kd tree work only on sparse matrices
final_X_tr=final_X_tr.todense()
final_X_test=final_X_test.todense()
final_X_cv=final_X_cv.todense()
```


In [55]:

```
#Calculating for finding Best K
#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
k_value=[]
roc_tr=[]
roc_cv=[]
max_auc_score=0
best_k=0
for i in range(1,50,2):
    # instantiate learning model (k = 50)
    knn = KNeighborsClassifier(n_neighbors=i,algorithm='kd_tree',metric='minkowski')

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

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

    # predict the response on the training
    pred_tr = knn.predict_proba(final_X_tr)
    pred_tr=(pred_tr)[:,1]
    roc_tr.append(roc_auc_score(y_tr,pred_tr))
    k_value.append(i)

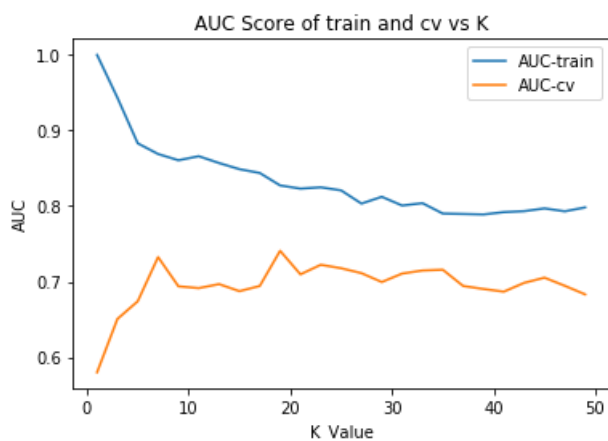
    #finding best k using loop
    if roc_auc_score(y_cv,pred_cv)>max_auc_score:
        best_k=i
        max_auc_score=roc_auc_score(y_cv,pred_cv)

print(best_k)
print(max_auc_score)
k_set5=best_k
auc_set5=max_auc_score
```

19
0.7407846715328467

In [56]:

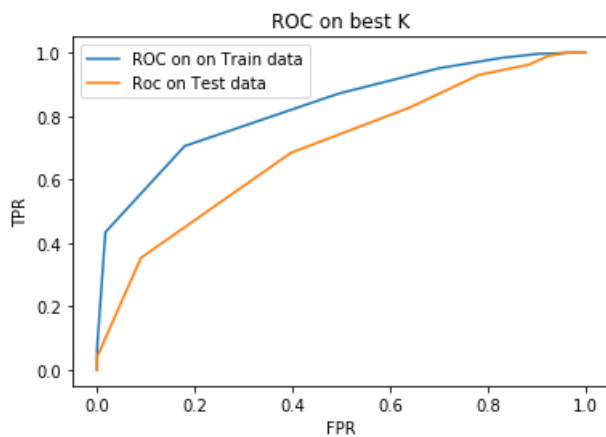
```
# plotting curve between AUC of cv and train with k
plt.plot(k_value,roc_tr,label="AUC-train")
plt.plot(k_value,roc_cv ,label="AUC-cv")
plt.legend()
plt.xlabel('K_Value')
plt.ylabel('AUC')
plt.title('AUC Score of train and cv vs K')
plt.show()
```



In [57]:

```
#1) Training the model using best K
knn = KNeighborsClassifier(n_neighbors=best_k, algorithm='kd_tree', metric='minkowski')
knn.fit(final_X_tr, y_tr)
#predicting probability of success Training data
pred_tr = knn.predict_proba(final_X_tr)
pred_tr=(pred_tr)[: ,1]
#predicting probability of success on Test data
pred_test = knn.predict_proba(final_X_test)
pred_test=(pred_test)[: ,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 K')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



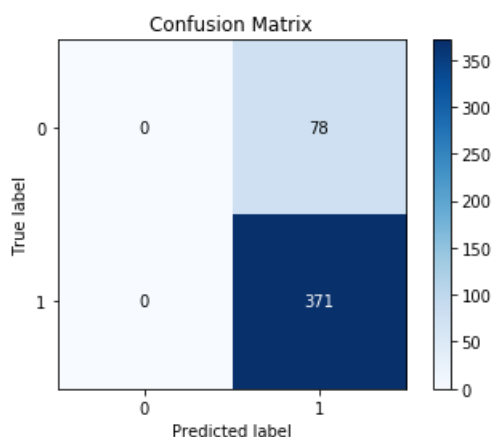
In [58]:

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

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

Out[58]:

<matplotlib.axes._subplots.AxesSubplot at 0x10bcb6d0>



[5.2.2] Applying KNN kd-tree on TFIDF, SET 6

In [59]:

```
#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)
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)

#converting them to dense bcz kd tree does not work on sparse matrices
final_X_tr=final_X_tr.todense()
final_X_test=final_X_test.todense()
final_X_cv=final_X_cv.todense()
```

In [99]:

```
#Calculating for finding Best K
#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
k_value=[]
roc_tr=[]
roc_cv=[]
max_auc_score=0
best_k=0
for i in range(1,50,2):
    # instantiate learning model (k = 50)
    knn = KNeighborsClassifier(n_neighbors=i, algorithm='kd_tree', metric='minkowski')

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

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

    # predict the response on the training
    pred_tr = knn.predict_proba(final_X_tr)
    pred_tr=(pred_tr)[:,1]
    roc_tr.append(roc_auc_score(y_tr, pred_tr))
    k_value.append(i)

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

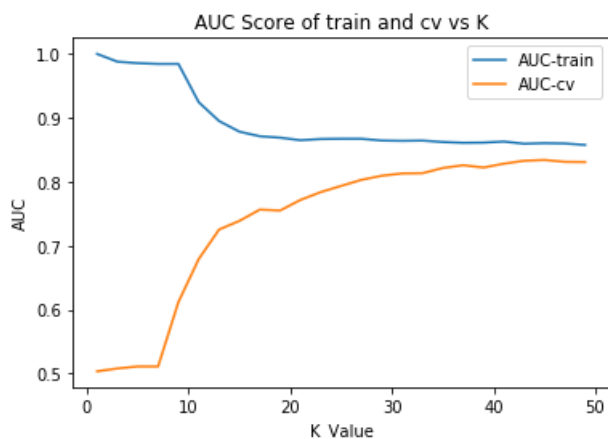
print(best_k)
print(max_auc_score)
k_set6=best_k
auc_set6=max_auc_score
```

```
41
0.6050993368198198
```

In [62]:

```
# plotting curve between AUC of cv and train with k
plt.plot(k_value, roc_tr, label="AUC-train")
plt.plot(k_value, roc_cv, label="AUC-cv")
```

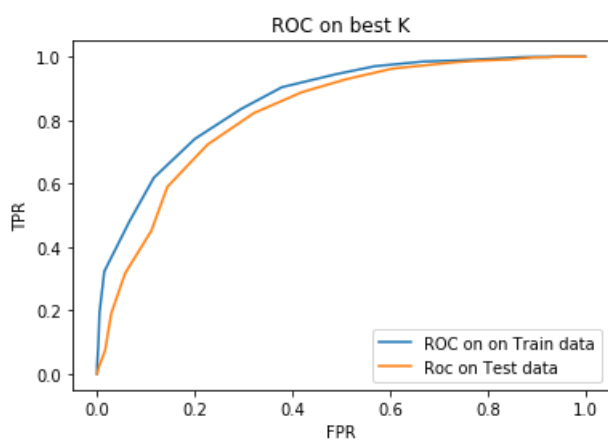
```
plt.legend()
plt.xlabel('K_Value')
plt.ylabel('AUC')
plt.title('AUC Score of train and cv vs K')
plt.show()
```



In [63]:

```
#1) Training the model using best K
knn = KNeighborsClassifier(n_neighbors=best_k, algorithm='kd_tree', metric='minkowski')
knn.fit(final_X_tr, y_tr)
#predicting probability of success Training data
pred_tr = knn.predict_proba(final_X_tr)
pred_tr=(pred_tr)[: ,1]
#predicting probability of success on Test data
pred_test = knn.predict_proba(final_X_test)
pred_test=(pred_test)[: ,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 K')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



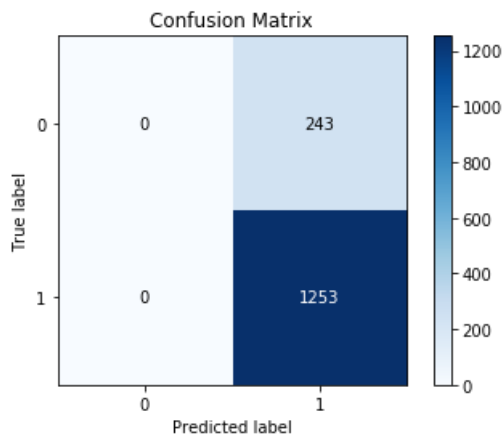
In [81]:

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

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

Out[81]:

<matplotlib.axes._subplots.AxesSubplot at 0x125294d0>



[5.2.3] Applying KNN kd-tree on AVG W2V, SET 7

In [82]:

```
#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()
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
    if cnt_words != 0:
        sent_vec /= cnt_words
    final_X_test.append(sent_vec)

```

```

100%|████████████████████████████████████████████████████████████████████████████████| 2443/2443 [00:04<00:00, 489.26it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 1047/1047 [00:02<00:00, 469.27it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 1496/1496 [00:03<00:00, 467.62it/s]

```

In [100]:

```

#Calculating for finding Best K
#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
k_value=[]
roc_tr=[]
roc_cv=[]
max_auc_score=0
best_k=0
for i in range(1,50,2):
    # instantiate learning model (k = 50)
    knn = KNeighborsClassifier(n_neighbors=i,algorithm='kd_tree',metric='minkowski')

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

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

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

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

print(best_k)
print(max_auc_score)
k_set7=best_k
auc_set7=max_auc_score

```

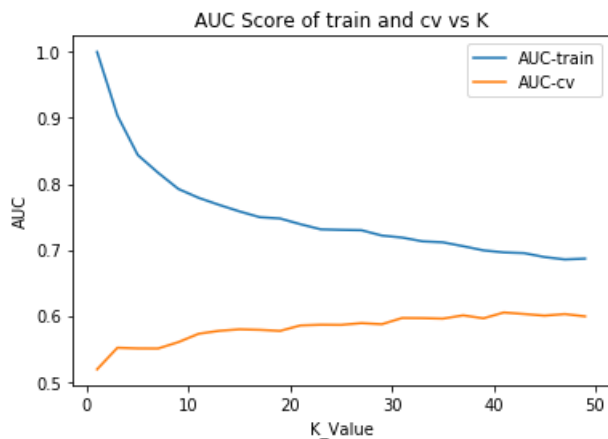
```

41
0.6050993368198198

```

In [89]:

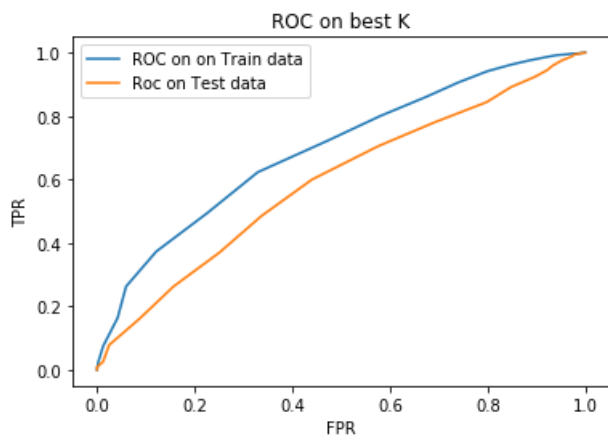
```
# plotting curve between between AUC of cv and train with k
plt.plot(k_value,roc_tr,label="AUC-train")
plt.plot(k_value,roc_cv ,label="AUC-cv")
plt.legend()
plt.xlabel('K_Value')
plt.ylabel('AUC')
plt.title('AUC Score of train and cv vs K')
plt.show()
```



In [90]:

```
#1)Training the model using best K
knn = KNeighborsClassifier(n_neighbors=best_k,algorithm='kd_tree',metric='minkowski')
knn.fit(final_X_tr, y_tr)
#predicting probability of success Training data
pred_tr = knn.predict_proba(final_X_tr)
pred_tr=(pred_tr)[:,1]
#predicting probability of success on Test data
pred_test = knn.predict_proba(final_X_test)
pred_test=(pred_test)[:,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 K')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



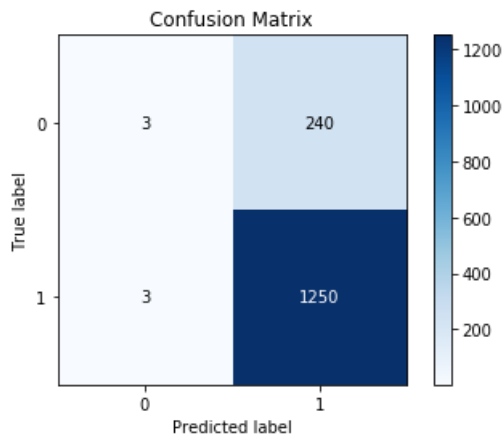
In [91]:

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

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

Out [91]:

<matplotlib.axes._subplots.AxesSubplot at 0x119fae50>



[5.2.4] Applying KNN kd-tree on TFIDF W2V, SET 8

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 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
        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 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)
```



```

        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]
            # 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
    final_X_test.append(sent_vec)
    row += 1

```

```

100%|████████████████████████████████████████| 2443/2443 [00:28<00:00, 89.31it/s]
100%|████████████████████████████████████████| 1047/1047 [00:11<00:00, 87.82it/s]
100%|████████████████████████████████████████| 1496/1496 [00:16<00:00, 92.23it/s]

```

In [101]:

```

#Calculating for finding Best K
#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
k_value=[]
roc_tr=[]
roc_cv=[]
max_auc_score=0
best_k=0
for i in range(1,50,2):
    # instantiate learning model (k = 50)
    knn = KNeighborsClassifier(n_neighbors=i,algorithm='kd_tree',metric='minkowski')

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

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

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

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

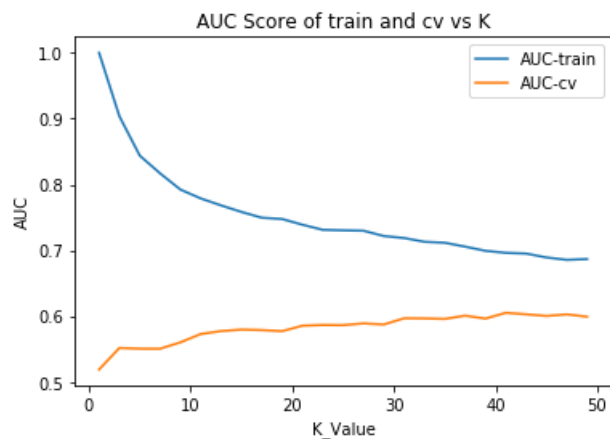
print(best_k)
print(max_auc_score)
k_set8=best_k
auc_set8=max_auc_score

```

41
0.6050993368198198

In [92]:

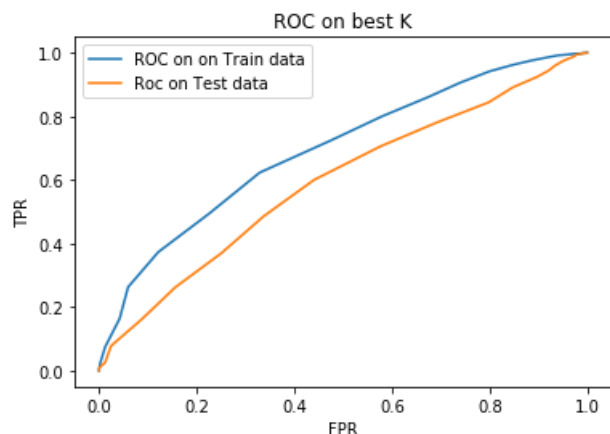
```
# plotting curve between between AUC of cv and train with k
plt.plot(k_value,roc_tr,label="AUC-train")
plt.plot(k_value,roc_cv ,label="AUC-cv")
plt.legend()
plt.xlabel('K_Value')
plt.ylabel('AUC')
plt.title('AUC Score of train and cv vs K')
plt.show()
```



In [93]:

```
#1)Training the model using best K
knn = KNeighborsClassifier(n_neighbors=best_k,algorithm='brute',metric='minkowski')
knn.fit(final_X_tr, y_tr)
#predicting probability of success Training data
pred_tr = knn.predict_proba(final_X_tr)
pred_tr=(pred_tr)[: ,1]
#predicting probability of success on Test data
pred_test = knn.predict_proba(final_X_test)
pred_test=(pred_test)[: ,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 K')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



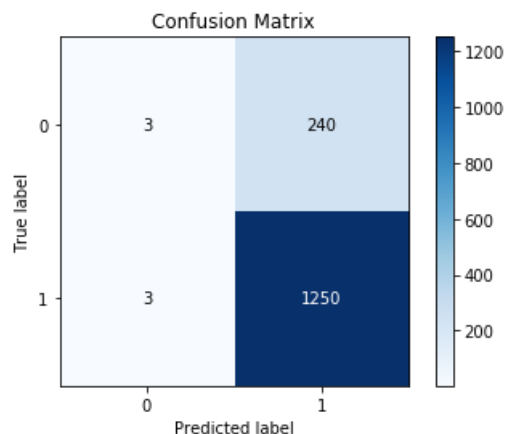
In [94]:

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

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

Out[94]:

<matplotlib.axes._subplots.AxesSubplot at 0x452abf0>



[6] Conclusions

In [104]:

```
# 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", "Brute", k_set1, auc_set1])
x.add_row(["TFIDF", "Brute", k_set2, auc_set2])
x.add_row(["AwgW2V", "Brute", k_set3, auc_set3])
x.add_row(["TFIDF-W2V", "Brute", k_set4, auc_set4])
x.add_row(["BOW", "k_d tree", k_set5, auc_set5])
x.add_row(["TFIDF", "k_d tree", k_set6, auc_set6])
x.add_row(["AwgW2V", "k_d tree", k_set7, auc_set7])
x.add_row(["TFIDF-W2V", "k_d tree", k_set8, auc_set8])
print(x)
```

Vectorizer	Model_name	Hyperamete-best_k	AUC
BOW	Brute	41	0.6050993368198198
TFIDF	Brute	45	0.8345484626277638
AwgW2V	Brute	41	0.6050993368198198
TFIDF-W2V	Brute	41	0.6050993368198198
BOW	k_d tree	19	0.7407846715328467
TFIDF	k_d tree	41	0.6050993368198198
AwgW2V	k_d tree	41	0.6050993368198198
TFIDF-W2V	k_d tree	41	0.6050993368198198