

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 the Score/Rating. A rating of 4 or 5 could be considered a positive review. A review of 1 or 2 could be considered negative. A review of 3 is neutral and ignored. This is an approximate and proxy way of determining the polarity (positivity/negativity) of a review.

Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

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

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

```
In [2]: %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:\Anaconda\lib\site-packages\gensim\utils.py:1209: UserWarning: detected
Windows; aliasing chunkize to chunkize_serial
  warnings.warn("detected Windows; aliasing chunkize to chunkize_seria
l")

```

[1]. Reading Data

```

In [3]: # using the SQLite Table to read data.
con = sqlite3.connect('database.sqlite')
#filtering only positive and negative reviews i.e.
# not taking into consideration those reviews with Score=3
# SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000, will give top 50
0000 data points
# you can change the number to any other number based on your computing
power

```

```
# filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE 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, and reviews with a score <3 a negative rating.
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[3]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator |
|---|----|------------|----------------|-------------|----------------------|------------------------|
| 0 | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1 | 1 |
| 1 | 2 | B00813GRG4 | A1D87F6ZCVE5NK | dll pa | 0 | 1 |

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenomin |
|---|----|------------|---------------|------------------------------------------|----------------------|--------------------|
| 2 | 3 | B000LQOCH0 | ABXLMWJIXXAIN | Natalia Corres "Natalia Corres" | 1 | |

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

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

(80668, 7)

Out[5]:

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

| | UserId | ProductId | ProfileName | Time | Score | Text | COUNT(*) |
|---|--------------------|------------|-----------------------|------------|-------|---------------------------------------------------|----------|
| 4 | #oc-R12KPBODL2B5ZD | B007OSBE1U | Christopher P. Presta | 1348617600 | 1 | I didnt like this coffee. Instead of telling y... | 2 |

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

Out[6]:

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

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

Out[7]: 393063

Exploratory Data Analysis

[2] Data Cleaning: Deduplication

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

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

Out[8]:

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

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

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

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

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

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

```
In [9]: #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 [10]: #Deduplication of entries
final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time",
"Text"}, keep='first', inplace=False)
final.shape
```

```
Out[10]: (4986, 10)
```

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

```
Out[11]: 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 [12]: display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND Id=44737 OR Id=64422
ORDER BY ProductID
""", con)
```



```
display.head()
```

Out[12]:

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

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

```
In [14]: #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[14]: 1    4178  
0      808  
Name: Score, dtype: int64
```

[3]. Text Preprocessing.

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

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

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

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

```
In [15]: # 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?<
br />http://www.amazon.com/VICTOR-FLY-MAGNET-BAIT-REFILL/dp/B00004RBDY<
br /><br />The Victor M380 and M502 traps are unreal, of course -- tota
l 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 Nabiso'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.

=====

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

=====

```
In [16]: # 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)
```

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]: *# 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.

I love to order my coffee on amazon. easy and shows up quickly. This k cup is great coffee. dcaf 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 it 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's 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 [18]: #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 [19]: #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 m sorry but these reviews do nobody any good beyond reminding us to look before ordering br br 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 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 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 br br 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

```
In [20]: # 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", 'it
```

```
s', 'itself', 'they', 'them', 'their', \
    'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'th
is', 'that', "that'll", 'these', 'those', \
    'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'h
ave', 'has', 'had', 'having', 'do', 'does', \
    'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or',
    'because', 'as', 'until', 'while', 'of', \
    'at', 'by', 'for', 'with', 'about', 'against', 'between',
    'into', 'through', 'during', 'before', 'after', \
    'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out',
    'on', 'off', 'over', 'under', 'again', 'further', \
    'then', 'once', 'here', 'there', 'when', 'where', 'why', 'h
ow', 'all', 'any', 'both', 'each', 'few', 'more', \
    'most', 'other', 'some', 'such', 'only', 'own', 'same', 's
o', 'than', 'too', 'very', \
    's', 't', 'can', 'will', 'just', 'don', "don't", 'should',
    "should've", 'now', 'd', 'll', 'm', 'o', 're', \
    've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't",
    'didn', "didn't", 'doesn', "doesn't", 'hadn', \
    "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "is
n't", 'ma', 'mightn', "mightn't", 'mustn', \
    "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn',
    "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
    'won', "won't", 'wouldn', "wouldn't"]])
```

```
In [23]: # Combining all the above students
from tqdm import tqdm
from bs4 import BeautifulSoup
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
```



```
100%|███████████| 4986/4986 [00:01<00:00, 2735.03it/s]
```

```
Out[24]: 'wow far two two star reviews one obviously no idea ordering wants cris
py cookies hey sorry reviews nobody good beyond reminding us look order
ing chocolate oatmeal cookies not like combination not order type cooki
e find combo quite nice really oatmeal sort calms rich chocolate flavor
gives cookie sort coconut type consistency let also remember tastes dif
fer given opinion soft chewy cookies advertised not crispy cookies blur
b would say crispy rather chewy happen like raw cookie dough however no
t see taste like raw cookie dough soft however confusion yes stick toge
ther soft cookies tend not individually wrapped would add cost oh yeah
chocolate chip cookies tend somewhat sweet want something hard crisp su
ggest nabisco ginger snaps want cookie soft chewy tastes like combinatio
n chocolate oatmeal give try place second order'
```

```
In [25]: ## Similarly you can do preprocessing for review summary also.
from tqdm import tqdm
preprocessed_summary = []
#tqdm is for printing the status bar
for sentence in tqdm(final['Summary'].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_summary.append(sentence.strip())
```

```
100%|███████████████████████████████████████████████████████████████████████████  
██████████| 4986/4986 [00:01<00:00, 3468.14it/s]
```

```
preprocessed_summary[1500]
```

'coffee supposedly premium tastes watery thin not good maybe old not sure waste using line bottom sitting shoes trash cans rained luggage also smells used not drink not buy'

[4] Featurization

[4.1] BAG OF WORDS

```
score = final['Score'] #storing all scores in new series
print(type(score))
print(score.shape)
```

```
<class 'pandas.core.series.Series'>
(4986,)
```

```
#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
```

```
In [29]: from sklearn.preprocessing import StandardScaler
std_data = StandardScaler(with_mean=False).fit_transform(final_counts)
std_data.shape
```

```
C:\Anaconda\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning: Data with input dtype int64 was converted to float64 by StandardScaler.
  warnings.warn(msg, DataConversionWarning)
C:\Anaconda\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning: Data with input dtype int64 was converted to float64 by StandardScaler.
  warnings.warn(msg, DataConversionWarning)
```

```
Out[29]: (4986, 12997)
```

```
In [30]: type(std_data)
```

```
Out[30]: scipy.sparse.csr.csr_matrix
```

```
In [31]: std_data = std_data.todense()
type(std_data)
```

```
Out[31]: numpy.matrixlib.defmatrix.matrix
```

[4.2] Bi-Grams and n-Grams.

```
In [32]: #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 [33]: 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 [34]: # 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 [35]: # Using Google News Word2Vectors

# in this project we are using a pretrained model by google
# its 3.3G file, once you load this into your memory
# it occupies ~9Gb, so please do this step only if you have >12G of ram
# we will provide a pickle file wich contains a dict ,
# and it contains all our courpus words as keys and model[word] as val
ues
# To use this code-snippet, download "GoogleNews-vectors-negative300.bi
n"
# from https://drive.google.com/file/d/0B7XkCwpI5KDYNlNUTTlSS21pQmM/edi
t
# it's 1.9GB in size.

# http://kavita-ganesan.com/gensim-word2vec-tutorial-starter-code/#.W17
SRFAzZPY
# you can comment this whole cell
# or change these 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 occured 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 gogole's word2vec file, keep want_to_train_w2v = True, to train your own w2v ")

[('snack', 0.9960731863975525), ('alternative', 0.9960423111915588),
 ('mellow', 0.9956516623497009), ('satisfying', 0.995613157749176), ('tasty', 0.9955406188964844), ('excellent', 0.9955399036407471), ('wonderful', 0.9952979683876038), ('bad', 0.9952245354652405), ('especially', 0.9951435923576355), ('crunch', 0.9951103925704956)]
=====
[('comes', 0.9995228052139282), ('actual', 0.9994649887084961), ('beef', 0.9994617700576782), ('oh', 0.9994484782218933), ('sound', 0.9994481205940247), ('wow', 0.9994404911994934), ('wife', 0.9994165897369385), ('major', 0.9994014501571655), ('person', 0.999396562576294), ('lays', 0.9993849992752075)]

```

```

In [36]: 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', 'love', 'call', 'instead', 'removed', 'easily', 'daughter', 'designed', 'printed', 'use', 'car', 'windows', '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 wAvg W2V, TFIDF-W2V

[4.4.1.1] Avg W2v

```
In [37]: # 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, yo
u might need to change this to 300 if you use google's w2v
    cnt_words = 0; # num of words with a valid vector in the sentence/re
view
    for word in sent: # for each word in a review/sentence
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    sent_vectors.append(sent_vec)
print(len(sent_vectors))
print(len(sent_vectors[0]))
```

```
100%|████████████████████████████████████████████████████████████████████████████████|
████████████████████████████████████████████████████████████████████████████████| 4986/4986 [00:03<00:00, 1247.48it/s]
```

```
4986
50
```

[4.4.1.2] TFIDF weighted W2v

```
In [38]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer()
model.fit(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 [39]: # 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)
    row += 1
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 4986/4986 [00:23<00:00, 212.18it/s]
```

[5] Applying TSNE

1. you need to plot 4 tsne plots with each of these feature set

- A. Review text, preprocessed one converted into vectors using (BOW)
 - B. Review text, preprocessed one converted into vectors using (TFIDF)
 - C. Review text, preprocessed one converted into vectors using (AVG W2v)
 - D. Review text, preprocessed one converted into vectors using (TFIDF W2v)
2. [Note 1: The TSNE accepts only dense matrices](#)
 3. [Note 2: Consider only 5k to 6k data points](#)

TSNE is an iterative algorithm: tries to process all the data in iterations.

Eventually it reaches a stage of clusters are not more moving.

Parameters of TSNE : perplexity, StepSize.

```
In [40]: # https://github.com/pavlin-polcar/fastTSNE you can try this also, this version is little faster than sklearn
import numpy as np
from sklearn.manifold import TSNE
from sklearn import datasets
import pandas as pd
import matplotlib.pyplot as plt

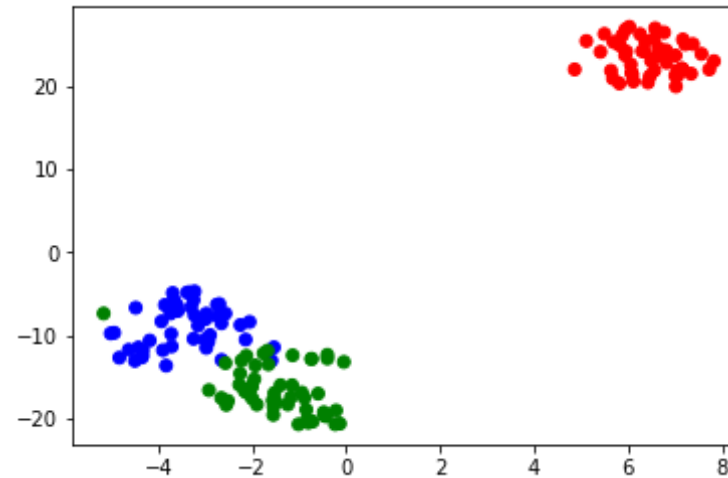
iris = datasets.load_iris()
x = iris['data']
y = iris['target']

tsne = TSNE(n_components=2, perplexity=30, learning_rate=200)

X_embedding = tsne.fit_transform(x)
# if x is a sparse matrix you need to pass it as X_embedding = tsne.fit_transform(x.toarray()) , .toarray() will convert the sparse matrix into dense matrix

for_tsne = np.hstack((X_embedding, y.reshape(-1,1)))
```

```
for_tsne_df = pd.DataFrame(data=for_tsne, columns=['Dimension_x', 'Dimension_y', 'Score'])
colors = {0: 'red', 1: 'blue', 2: 'green'}
plt.scatter(for_tsne_df['Dimension_x'], for_tsne_df['Dimension_y'], c=for_tsne_df['Score'].apply(lambda x: colors[x]))
plt.show()
```



[5.1] Applying TNSE on Text BOW vectors

```
In [40]: # please write all the code with proper documentation, and proper titles for each subsection
# when you plot any graph make sure you use
    # a. Title, that describes your plot, this will be very helpful to the reader
    # b. Legends if needed
    # c. X-axis label
    # d. Y-axis label

from sklearn.preprocessing import StandardScaler
print(final_bigram_counts.shape)
std_data = StandardScaler(with_mean=False).fit_transform(final_bigram_counts)
```

```
print(std_data.shape)
type(std_data)
std_data=std_data.todense()
print(type(std_data))
```

(4986, 3144)

C:\Anaconda\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning: Data with input dtype int64 was converted to float64 by StandardScaler.

warnings.warn(msg, DataConversionWarning)

C:\Anaconda\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning: Data with input dtype int64 was converted to float64 by StandardScaler.

warnings.warn(msg, DataConversionWarning)

(4986, 3144)

<class 'numpy.matrixlib.defmatrix.matrix'>

Perplexity: No of neighbors loosely to when the distances preserved.

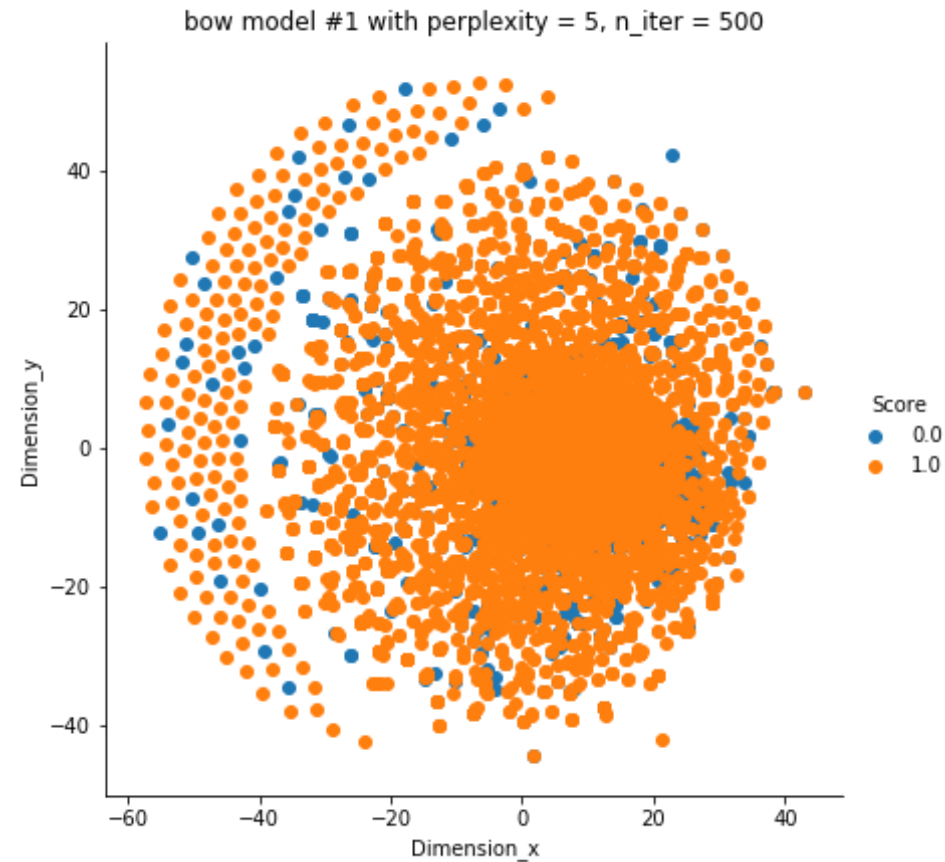
StepSize: No of iterations better the solution.

```
In [41]: from sklearn.manifold import TSNE
from sklearn import datasets
import pandas as pd
import matplotlib.pyplot as plt

model = TSNE(n_components=2, perplexity=5, n_iter=500, learning_rate=200, random_state=0)
for_tsne = model.fit_transform(std_data)
# if x is a sparse matrix you need to pass it as X_embedding = tsne.fit_transform(x.toarray()) , .toarray() will convert the sparse matrix into dense matrix

for_tsne = np.vstack((for_tsne.T, score)).T
```

```
for_tsne_df = pd.DataFrame(data=for_tsne, columns=['Dimension_x', 'Dimension_y', 'Score'])
sns.FacetGrid(for_tsne_df, hue="Score", height=6).map(plt.scatter, 'Dimension_x', 'Dimension_y').add_legend()
plt.title('bow model #1 with perplexity = 5, n_iter = 500')
plt.show()
```



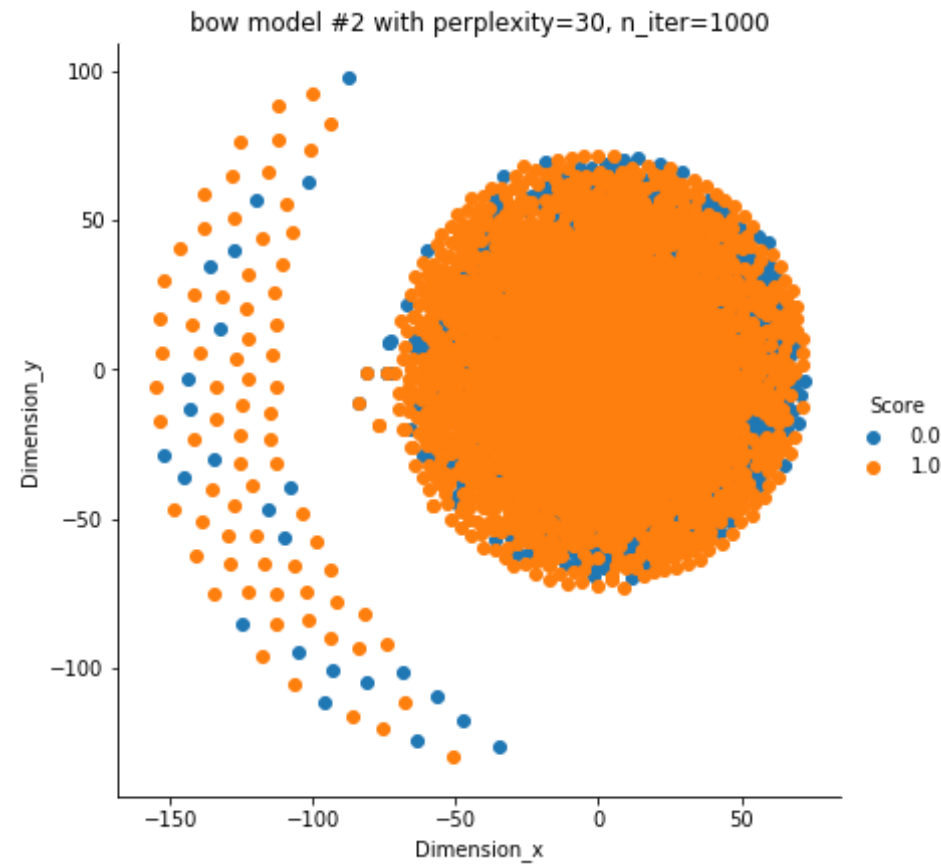
Observation:

1. Bow model at 500 iterations is not classifying positive and negative reviews of the preprocessed text it's unable to distinguish since the points formed cluster of overlapped(crowded) .

2. Let's increase the `n_iter`, perplexity and see whether it's going to distinguish b/w positive and negative text.

```
In [42]: model = TSNE(n_components=2, perplexity=30, n_iter=1000, learning_rate=
400, random_state=0)
for_tsne = model.fit_transform(std_data)
# if x is a sparse matrix you need to pass it as X_embedding = tsne.fit_
_transform(x.toarray()) , .toarray() will convert the sparse matrix int
o dense matrix

for_tsne = np.vstack((for_tsne.T, score)).T
for_tsne_df = pd.DataFrame(data=for_tsne, columns=['Dimension_x', 'Dimen
sion_y', 'Score'])
sns.FacetGrid(for_tsne_df, hue="Score", height=6).map(plt.scatter, 'Dim
ension_x', 'Dimension_y').add_legend()
plt.title('bow model #2 with perplexity=30, n_iter=1000')
plt.show()
```



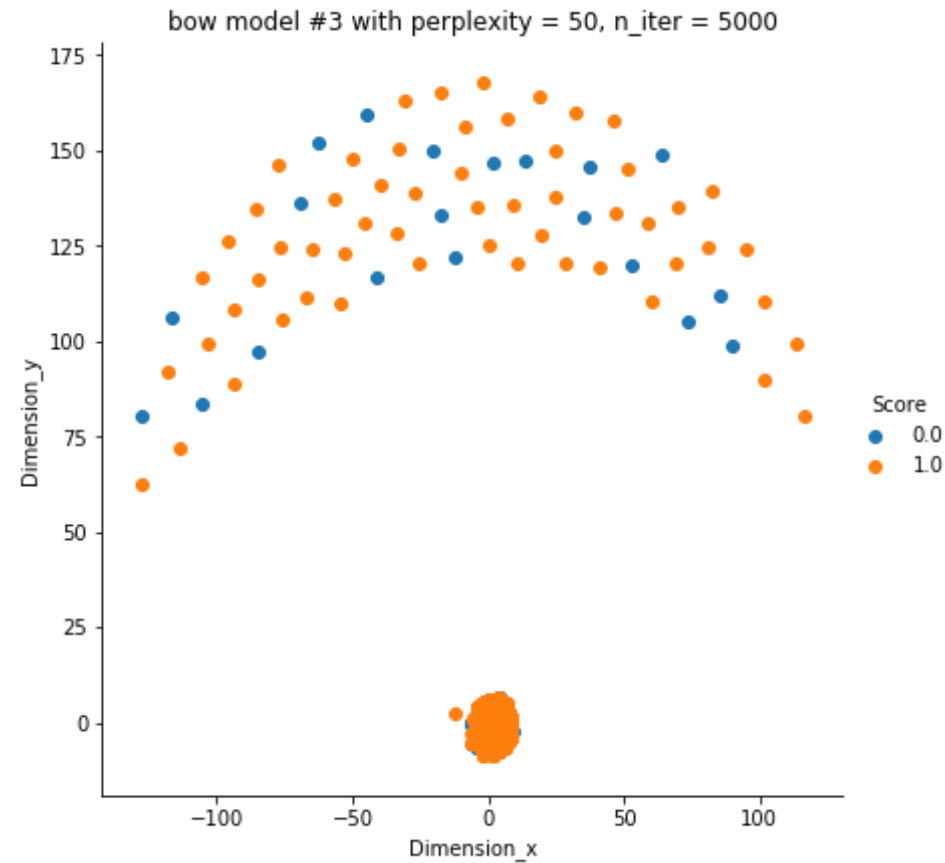
Observation:

1. Bow model2 with increment in iteration to 1000 is also unable to distinguish b/w positive and negative text.
2. let's increment the iteration and perplexity and check whether model is going to stable.
3. The points are more grouped as compared to above plot where points are scattered.

```
In [44]: model = TSNE(n_components=2, perplexity=50, n_iter=5000, learning_rate=
        600, random_state=0)
        for_tsne = model.fit_transform(std_data)
        # if x is a sparse matrix you need to pass it as X_embedding = tsne.fit
```

```
_transform(x.toarray()) , .toarray() will convert the sparse matrix into dense matrix
```

```
for_tsne = np.vstack((for_tsne.T, score)).T  
for_tsne_df = pd.DataFrame(data=for_tsne, columns=['Dimension_x', 'Dimension_y', 'Score'])  
sns.FacetGrid(for_tsne_df, hue="Score", height=6).map(plt.scatter, 'Dimension_x', 'Dimension_y').add_legend()  
plt.title('bow model #3 with perplexity = 50, n_iter = 5000')  
plt.show()
```



Observation:

1. Bow model3 with increased iteration of 5000, we can infer that model became stable. it's unable to classify text reviews.
2. Still we cannot determine whether the reviews of text are positive or not.

[5.1] Applying TNSE on Text TFIDF vectors

```
In [45]: # please write all the code with proper documentation, and proper titles for each subsection  
# when you plot any graph make sure you use  
# a. Title, that describes your plot, this will be very helpful to the reader  
# b. Legends if needed  
# c. X-axis label  
# d. Y-axis label
```

```
from sklearn.preprocessing import StandardScaler  
std_data1 = StandardScaler(with_mean = False).fit_transform(final_tf_idf)  
print(std_data1.shape)  
type(std_data1)  
std_data1 = std_data1.todense()  
type(std_data1)
```

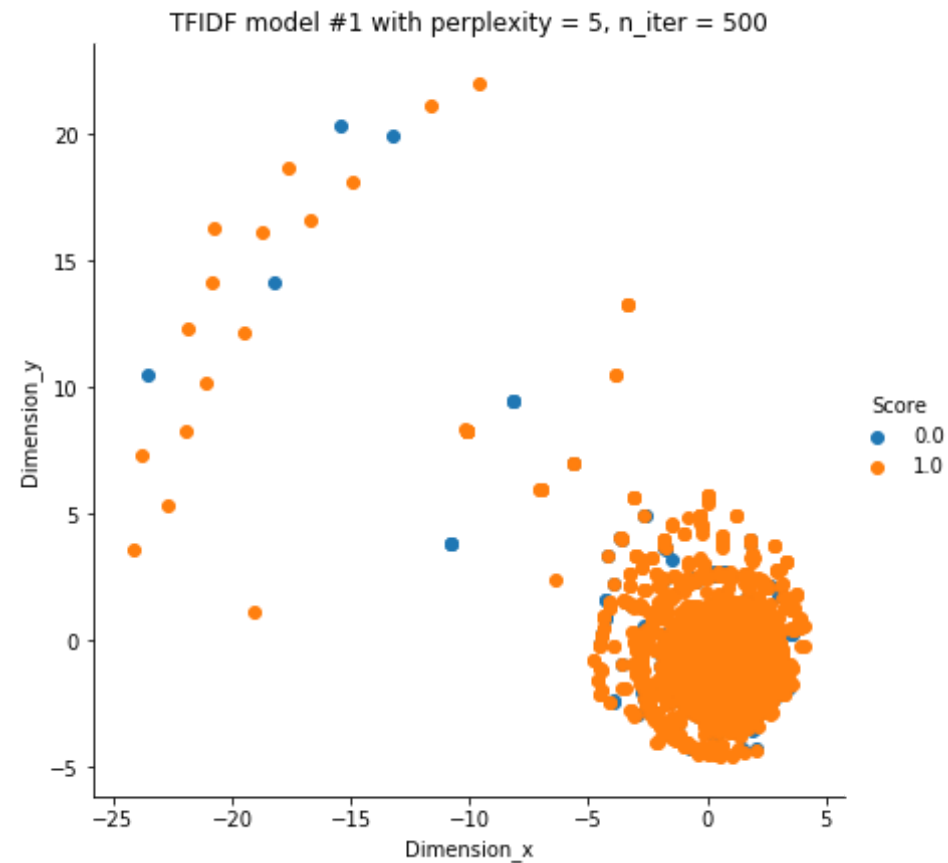
```
(4986, 3144)
```

```
Out[45]: numpy.matrixlib.defmatrix.matrix
```

```
In [57]: model = TSNE(n_components=2, perplexity=5, n_iter=500, learning_rate=200, random_state=0)  
for_tsne = model.fit_transform(std_data1)  
# if x is a sparse matrix you need to pass it as X_embedding = tsne.fit_transform(x.toarray()) , .toarray() will convert the sparse matrix into dense matrix  
  
for_tsne = np.vstack((for_tsne.T, score)).T  
for_tsne_df = pd.DataFrame(data=for_tsne, columns=['Dimension_x', 'Dimension_y'])
```



```
sion_y', 'Score'])  
sns.FacetGrid(for_tsne_df, hue="Score", height=6).map(plt.scatter, 'Dim  
ension_x', 'Dimension_y').add_legend()  
plt.title('TFIDF model #1 with perplexity = 5, n_iter = 500')  
plt.show()
```



Observation:

1. Plot became dispersed still it's not stable with iter(500).

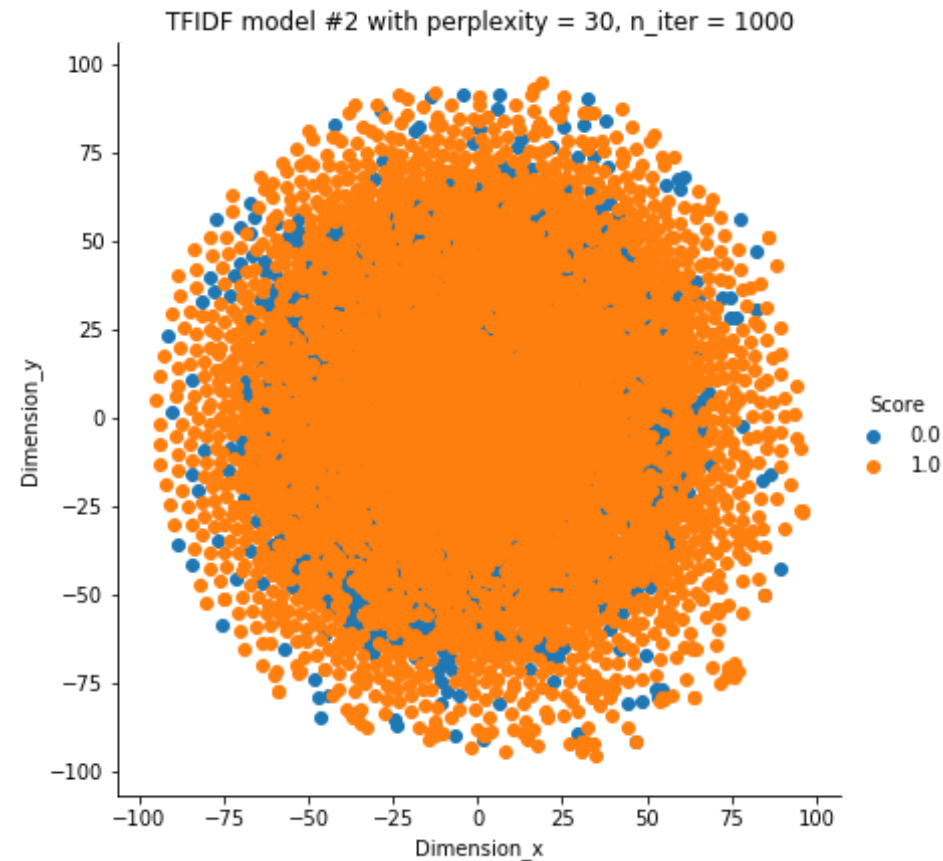
```
In [58]: model = TSNE(n_components=2, perplexity=30, n_iter=1000, learning_rate=
```

```

400, random_state=0)
for_tsne = model.fit_transform(std_data1)
# if x is a sparse matrix you need to pass it as X_embedding = tsne.fit
_transform(x.toarray()) , .toarray() will convert the sparse matrix into
a dense matrix

for_tsne = np.vstack((for_tsne.T, score)).T
for_tsne_df = pd.DataFrame(data=for_tsne, columns=['Dimension_x', 'Dimension_y', 'Score'])
sns.FacetGrid(for_tsne_df, hue="Score", height=6).map(plt.scatter, 'Dimension_x', 'Dimension_y').add_legend()
plt.title('TFIDF model #2 with perplexity = 30, n_iter = 1000')
plt.show()

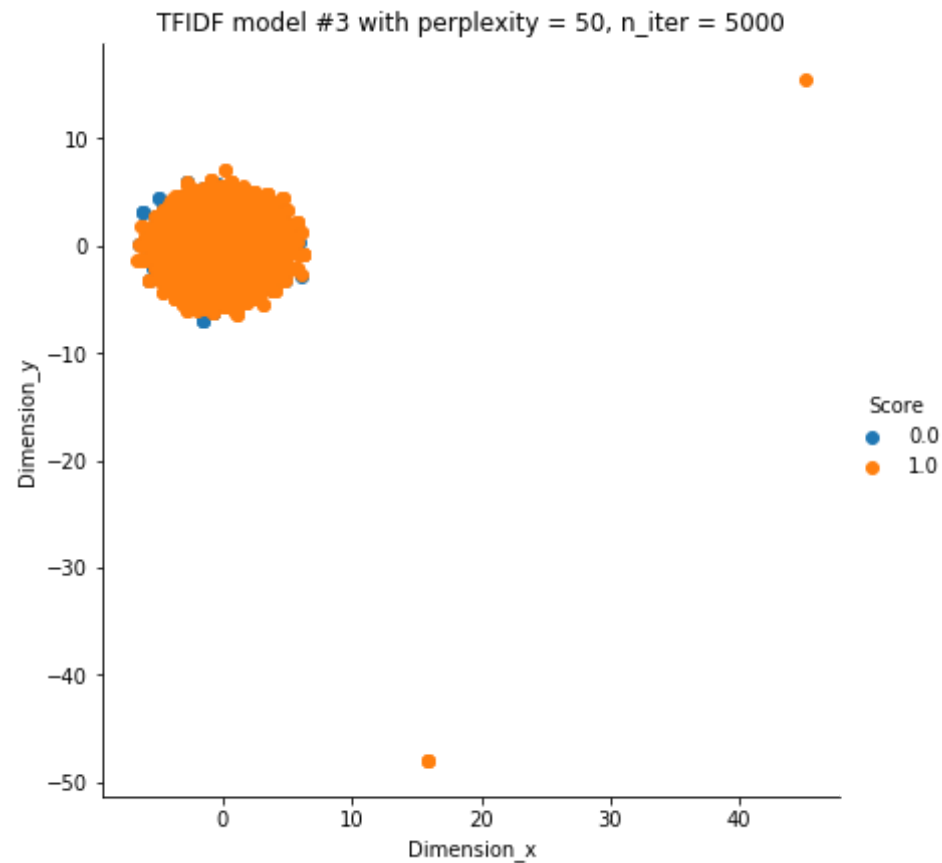
```



Observation: Plot became Stable at 1000 iter, due to class imbalance of more positive points and less negative points. '1' : Positive points. '2' : Negative Points.

```
In [59]: model = TSNE(n_components=2, perplexity=50, n_iter=5000, learning_rate=
600, random_state=0)
for_tsne = model.fit_transform(std_data1)
# if x is a sparse matrix you need to pass it as X_embedding = tsne.fit
_transform(x.toarray()) , .toarray() will convert the sparse matrix into
a dense matrix

for_tsne = np.vstack((for_tsne.T, score)).T
for_tsne_df = pd.DataFrame(data=for_tsne, columns=['Dimension_x', 'Dimension_y', 'Score'])
sns.FacetGrid(for_tsne_df, hue="Score", height=6).map(plt.scatter, 'Dimension_x', 'Dimension_y').add_legend()
plt.title('TFIDF model #3 with perplexity = 50, n_iter = 5000')
plt.show()
```



Observation: Most of the points are clustered around top left handside of the plot and is not stable.

[5.3] Applying TNSE on Text Avg W2V vectors

```
In [49]: # please write all the code with proper documentation, and proper titles for each subsection
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
```

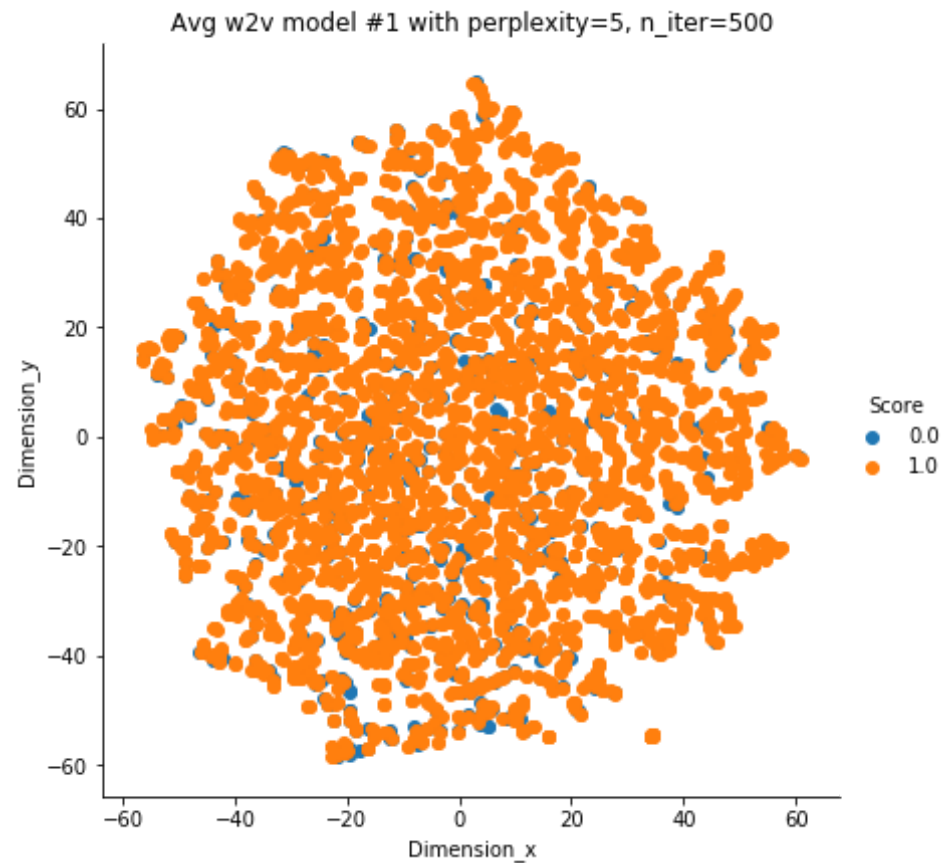
```
# b. Legends if needed  
# c. X-axis label  
# d. Y-axis label
```

```
from sklearn.preprocessing import StandardScaler  
std_data2 = StandardScaler(with_mean= False).fit_transform(sent_vectors  
)  
print(std_data2.shape)  
type(std_data2)
```

```
(4986, 50)
```

Out[49]: numpy.ndarray

```
In [50]: model = TSNE(n_components=2, perplexity=5, learning_rate=200, n_iter=50  
0, random_state=0)  
for_tsne = model.fit_transform(std_data2)  
# if x is a sparse matrix you need to pass it as X_embedding = tsne.fit  
_transform(x.toarray()) , .toarray() will convert the sparse matrix int  
o dense matrix  
  
for_tsne = np.vstack((for_tsne.T, score)).T  
for_tsne_df = pd.DataFrame(data=for_tsne, columns= ['Dimension_x', 'Dim  
ension_y', 'Score'])  
sns.FacetGrid(for_tsne_df, hue="Score", height=6).map(plt.scatter, 'Dim  
ension_x', 'Dimension_y').add_legend()  
plt.title('Avg w2v model #1 with perplexity=5, n_iter=500')  
plt.show()
```

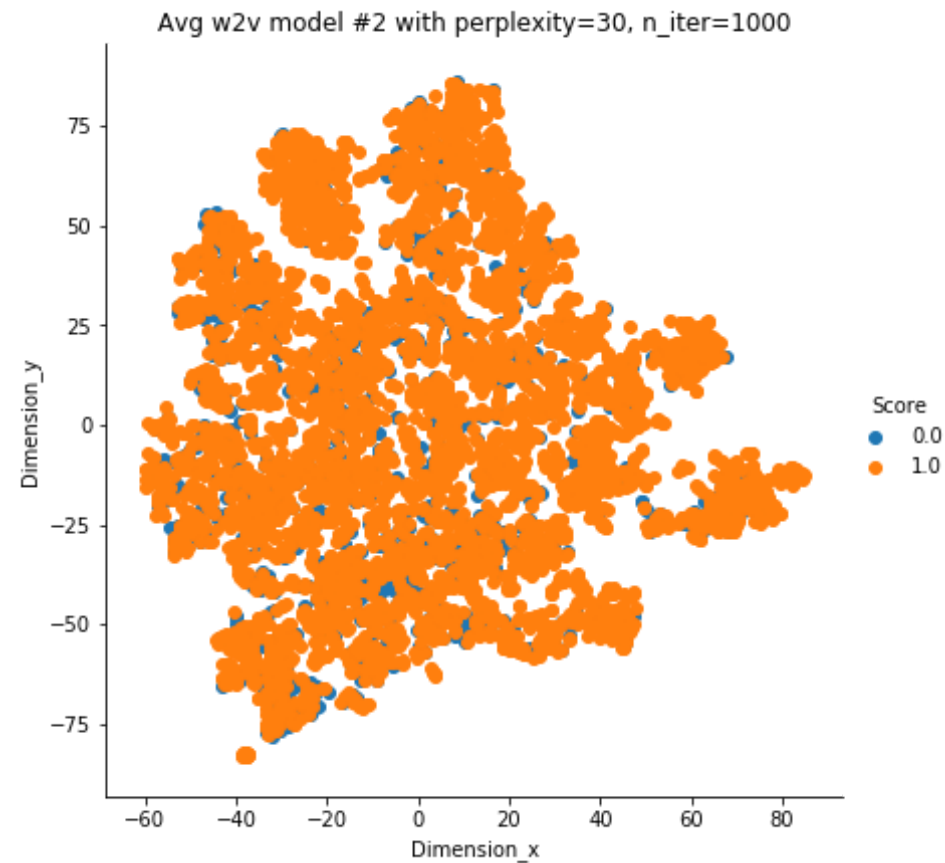


Observation:

1. Due to class imbalance plot is scattered at 500 iter. More no of Positive points are observed than negative points.

```
In [51]: model = TSNE(n_components=2, perplexity=30, learning_rate=400, n_iter=1000, random_state=0)
for_tsne = model.fit_transform(std_data2)
# if x is a sparse matrix you need to pass it as X_embedding = tsne.fit_transform(x.toarray()) , .toarray() will convert the sparse matrix into dense matrix
```

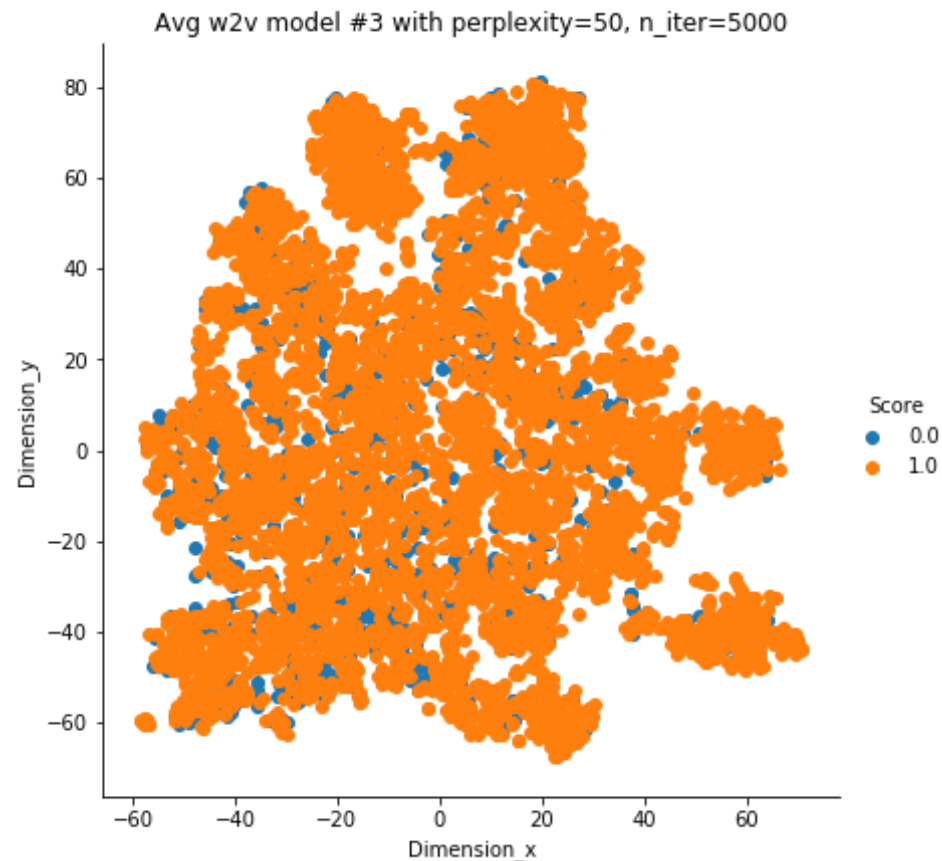
```
for_tsne = np.vstack((for_tsne.T, score)).T
for_tsne_df = pd.DataFrame(data=for_tsne, columns=['Dimension_x', 'Dimension_y', 'Score'])
sns.FacetGrid(for_tsne_df, hue="Score", height=6).map(plt.scatter, 'Dimension_x', 'Dimension_y').add_legend()
plt.title('Avg w2v model #2 with perplexity=30, n_iter=1000')
plt.show()
```



Observation:

1. we can observe that points overlap massively and form different clusters with more no of positive points.

```
In [52]: model = TSNE(n_components=2, perplexity=50, learning_rate= 600, n_iter=
5000, random_state=0)
for_tsne = model.fit_transform(std_data2)
for_tsne = np.vstack((for_tsne.T, score)).T
for_tsne_df = pd.DataFrame(data=for_tsne, columns=['Dimension_x', 'Dimension_y', 'Score'])
sns.FacetGrid(for_tsne_df, hue='Score', height=6).map(plt.scatter, 'Dimension_x', 'Dimension_y').add_legend()
plt.title('Avg w2v model #3 with perplexity=50, n_iter=5000')
plt.show()
```

Observation: At model 3 of Avg w2v with 5000 iter also we can see overlap of +ve and -ve points.

[5.4] Applying TNSE on Text TFIDF weighted W2V vectors

```
In [53]: # please write all the code with proper documentation, and proper titles for each subsection
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
```

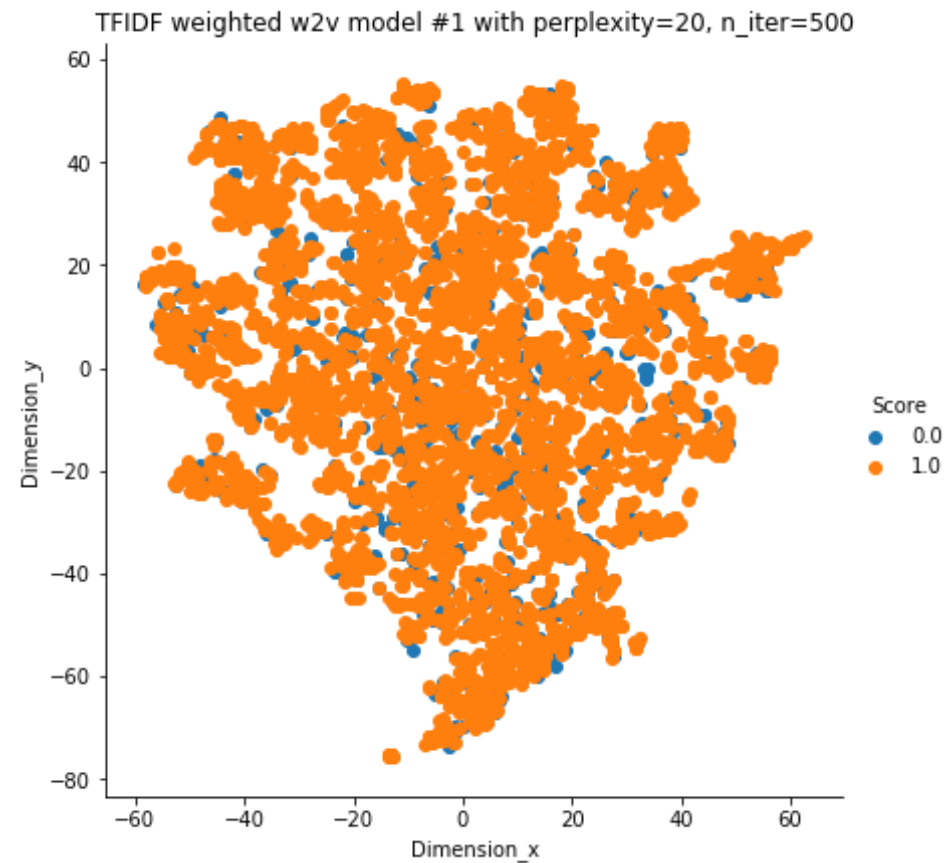
```
# b. Legends if needed  
# c. X-axis label  
# d. Y-axis label
```

```
from sklearn.preprocessing import StandardScaler  
std_data3 = StandardScaler(with_mean = False).fit_transform(tfidf_sent_  
vectors)  
print(std_data3.shape)  
type(std_data3)
```

```
(4986, 50)
```

```
Out[53]: numpy.ndarray
```

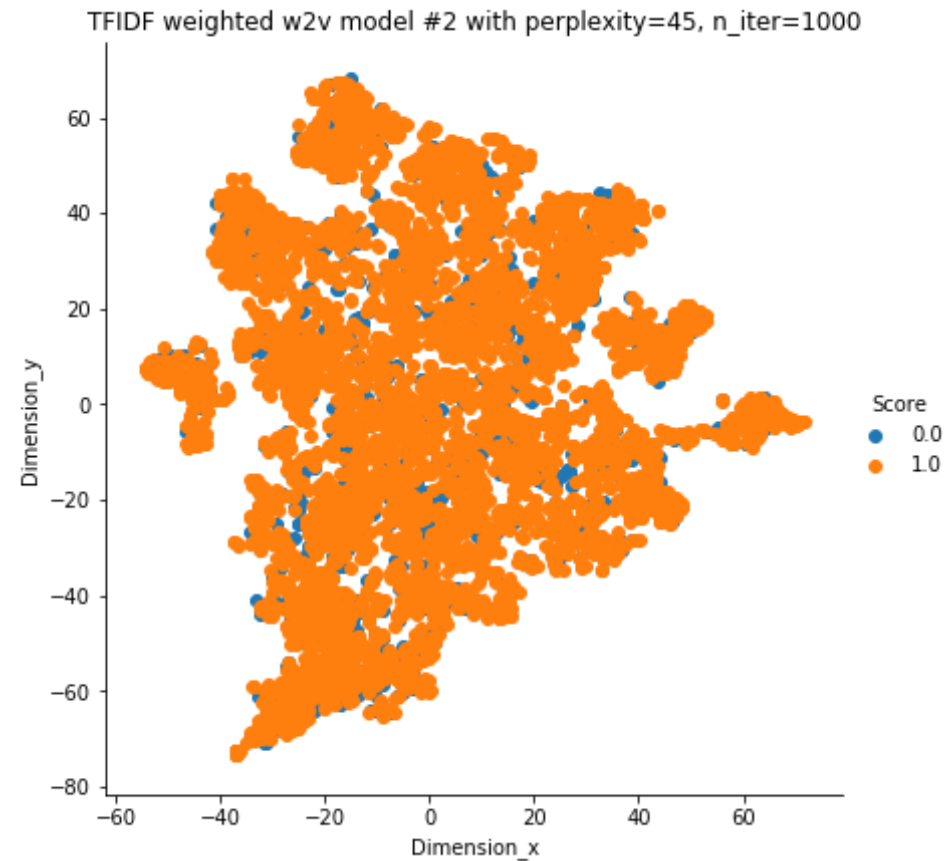
```
In [54]: model = TSNE(n_components=2, perplexity=20, n_iter=500, learning_rate=6  
00, random_state=0)  
for_tsne = model.fit_transform(std_data3)  
for_tsne = np.vstack((for_tsne.T, score)).T  
for_tsne_df = pd.DataFrame(data=for_tsne, columns=['Dimension_x', 'Dim  
ension_y', 'Score'])  
sns.FacetGrid(for_tsne_df, hue='Score', height=6).map(plt.scatter, 'Dim  
ension_x', 'Dimension_y').add_legend()  
plt.title('TFIDF weighted w2v model #1 with perplexity=20, n_iter=500')  
plt.show()
```



Observation: we can observe that both the classes are not going to be separable.

```
In [55]: model = TSNE(n_components=2, perplexity=45, n_iter=1000, learning_rate=
400, random_state=0)
for_tsne = model.fit_transform(std_data3)
for_tsne = np.vstack((for_tsne.T, score)).T
for_tsne_df = pd.DataFrame(data=for_tsne, columns=['Dimension_x', 'Dimension_y', 'Score'])
sns.FacetGrid(for_tsne_df, hue='Score', height=6).map(plt.scatter, 'Dimension_x', 'Dimension_y').add_legend()
plt.title('TFIDF weighted w2v model #2 with perplexity=45, n_iter=1000')
```

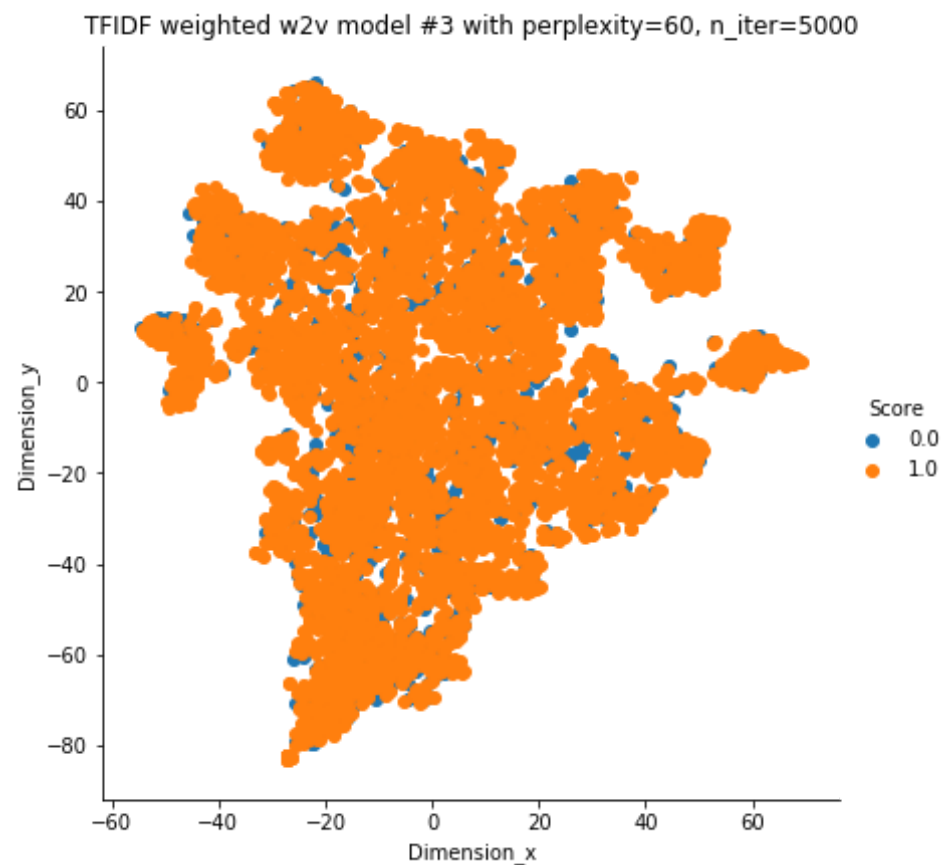
```
)  
plt.show()
```



Observation: None of the measurements of TSNE helpful to determine the class labels.

```
In [56]: model = TSNE(n_components=2, perplexity= 60, n_iter=5000, learning_rate  
=200, random_state=0)  
for_tsne = model.fit_transform(std_data3)  
for_tsne = np.vstack((for_tsne.T, score)).T  
for_tsne_df = pd.DataFrame(data=for_tsne, columns=['Dimension_x', 'Dimension_y', 'Score'])  
sns.FacetGrid(for_tsne_df, hue='Score', height=6).map(plt.scatter, 'Dim
```

```
ension_x', 'Dimension_y').add_legend()  
plt.title('TFIDF weighted w2v model #3 with perplexity=60, n_iter=5000'  
)  
plt.show()
```



Observation: By expermenting with increasing iterartions and perplexity also plots are not differntaited well.

[6] Conclusions

1. Experimenting with multiple no of iterations, and perplexity until the plot becomes stable.
2. Due to class imbalance we cannot determine the review points, as it belongs to positive class or negative class label. Certain words occur in both the class labels.
3. More no iterations and perplexity better is the plot. It becomes stable at certain iteration.
4. TSNE expands dense clusters, shrinks sparse clusters.
5. TSNE group points based on their visual similarity.
6. Points which are visually similar group togetehr, points less similar are far away.

In []: