# 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 - unqiue 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 cosnidered a positive review. A review of 1 or 2 could be considered negative. A review of 3 is nuetral 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 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 id above 3, then the recommendation wil 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
```

```python
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
```

# [1]. Reading Data

```python
In [2]: # 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 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, 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[2]:

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

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

In [55]:
```python
display = pd.read_sql_query("""
SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
FROM Reviews
GROUP BY UserId
HAVING COUNT(*)>1
""", con)
```

In [56]:
```python
print(display.shape)
display.head()
```

(80668, 7)

Out[56]:

| | UserId | ProductId | ProfileName | Time | Score | Text | COU |
|---|--------|-----------|-------------|------|-------|------|-----|
| 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 |

| | UserId | ProductId | ProfileName | Time | Score | Text | COU |
|---|---|---|---|---|---|---|---|
| 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 [57]: `display[display['UserId']=='AZY10LLTJ71NX']`

Out[57]:

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

In [58]: `display['COUNT(*)'].sum()`

Out[58]: 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 [3]:
```
display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID
""", con)
display.head()
```

Out[3]:

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

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfuln |
|---|---|---|---|---|---|---|
| 3 | 73791 | B000HDOPZG | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |
| 4 | 155049 | B000PAQ75C | AR5J8UI46CURR | Geetha Krishnan | 2 | 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 delelte 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.

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

In [5]:
```
#Deduplication of entries
final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time"
,"Text"}, keep='first', inplace=False)
final.shape
```

Out[5]: (4986, 10)

In [6]:
```
#Checking to see how much % of data still remains
(final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

Out[6]: 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 calcualtions

In [63]:
```
display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND Id=44737 OR Id=64422
ORDER BY ProductID
""", con)

display.head()
```

Out[63]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfuln |
|---|---|---|---|---|---|---|

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfuln |
|---|---|---|---|---|---|---|
| 0 | 64422 | B000MIDROQ | A161DK06JJMCYF | J. E. Stephens "Jeanne" | 3 | 1 |
| 1 | 44737 | B001EQ55RW | A2V0I904FH7ABY | Ram | 3 | 2 |

In [7]:
```
final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

In [8]:
```
#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[8]:
```
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 obsereved to be better than Porter Stemming)

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

```python
In [9]:  # 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 t
hese chips are.  The best thing was that there were a lot of "brown" ch
ips 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 fl
avor because they do not seem to be as salty, and the onion flavor is b
etter.  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 the
y were ordering; the other wants crispy cookies.  Hey, I'm sorry; but t
hese reviews do nobody any good beyond reminding us to look  before ord
ering.<br /><br />These are chocolate-oatmeal cookies.  If you don't li
ke that combination, don't order this type of cookie.  I find the combo
quite nice, really.  The oatmeal sort of "calms" the rich chocolate fla
vor 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 "c
rispy" cookies, or the blurb would say "crispy," rather than "chewy."
I happen to like raw cookie dough; however, I don't see where these tas
te like raw cookie dough.  Both are soft, however, so is this the confu
sion?  And, yes, they stick together.  Soft cookies tend to do that.  T
hey 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 Nabiso's Ginger Snaps.  If you
want a cookie that's soft, chewy and tastes like a combination of choco
late 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.<br />Thi
s k cup is great coffee.  dcaf is very good as well
=======================================================

In [10]:
```python
# remove urls from text python: https://stackoverflow.com/a/40823105/40
84039
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?<
br /> /><br />The Victor M380 and M502 traps are unreal, of course -- t
otal fly genocide. Pretty stinky, but only right nearby.

In [11]:
```
# 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 gen
ocide. Pretty stinky, but only right nearby.
==================================================
I recently tried this flavor/brand and was surprised at how delicious t
hese chips are.  The best thing was that there were a lot of "brown" ch
ips 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 fl avor because they do not seem to be as salty, and the onion flavor is b etter.  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 the y were ordering; the other wants crispy cookies.  Hey, I'm sorry; but t hese reviews do nobody any good beyond reminding us to look  before ord ering.These are chocolate-oatmeal cookies.  If you don't like that comb ination, don't order this type of cookie.  I find the combo quite nice, really.  The oatmeal sort of "calms" the rich chocolate flavor and give s the cookie sort of a coconut-type consistency.  Now let's also rememb er 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 co okie dough; however, I don't see where these taste like raw cookie doug h.  Both are soft, however, so is this the confusion?  And, yes, they s tick 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 s uggest 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 tr y.  I'm here to place my second order.
====================================================
love to order my coffee on amazon.  easy and shows up quickly.This k cu p is great coffee.  dcaf is very good as well

In [12]:
```python
# 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
```

```python
        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 [13]:
```python
sent_1500 = decontracted(sent_1500)
print(sent_1500)
print("="*50)
```

```
Wow.  So far, two two-star reviews.  One obviously had no idea what the
y were ordering; the other wants crispy cookies.  Hey, I am sorry; but
these reviews do nobody any good beyond reminding us to look  before or
dering.<br /><br />These are chocolate-oatmeal cookies.  If you do not
like that combination, do not order this type of cookie.  I find the co
mbo quite nice, really.  The oatmeal sort of "calms" the rich chocolate
flavor and gives the cookie sort of a coconut-type consistency.  Now le
t 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 "che
wy."  I happen to like raw cookie dough; however, I do not see where th
ese taste like raw cookie dough.  Both are soft, however, so is this th
e confusion?  And, yes, they stick together.  Soft cookies tend to do t
hat.  They are not individually wrapped, which would add to the cost.
Oh yeah, chocolate chip cookies tend to be somewhat sweet.<br /><br />S
o, if you want something hard and crisp, I suggest Nabiso is Ginger Sna
ps.  If you want a cookie that is soft, chewy and tastes like a combina
tion of chocolate and oatmeal, give these a try.  I am here to place my
second order.
==================================================
```

In [14]:
```python
#remove words with numbers python: https://stackoverflow.com/a/1808237
0/4084039
sent_0 = re.sub("\S*\d\S*", "", sent_0).strip()
print(sent_0)
```

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

In [15]:
```python
#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 Nabiso 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 [16]:
```python
# 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 [17]:  # Combining all the above stundents
          from tqdm import tqdm
          preprocessed_reviews = []
          # tqdm is for printing the status bar
          for sentance in tqdm(final['Text'].values):
              sentance = re.sub(r"http\S+", "", sentance)
              sentance = BeautifulSoup(sentance, 'lxml').get_text()
              sentance = decontracted(sentance)
              sentance = re.sub("\S*\d\S*", "", sentance).strip()
              sentance = re.sub('[^A-Za-z]+', ' ', sentance)
              # https://gist.github.com/sebleier/554280
              sentance = ' '.join(e.lower() for e in sentance.split() if e.lower
```

```
        () not in stopwords)
        preprocessed_reviews.append(sentance.strip())
```

100%|████████████████████████████████████| 4986/4986 [00:07<00:00, 64
8.43it/s]

In [18]:
```python
#adding a column of CleanedText which displays the data after pre-proce
ssing of the review
final['CleanedText']=np.array(preprocessed_reviews)
final['CleanedText']=final['CleanedText'].str.decode("utf-8")
#below the processed review can be seen in the CleanedText Column
print('Shape of final',final.shape)
final.head()
```

Shape of final (4986, 11)

Out[18]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfu |
|---|---|---|---|---|---|---|
| 2546 | 2774 | B00002NCJC | A196AJHU9EASJN | Alex Chaffee | 0 | 0 |
| 2547 | 2775 | B00002NCJC | A13RRPGE79XFFH | reader48 | 0 | 0 |
| 1145 | 1244 | B00002Z754 | A3B8RCEI0FXFI6 | B G Chase | 10 | 10 |
| 1146 | 1245 | B00002Z754 | A29Z5PI9BW2PU3 | Robbie | 7 | 7 |

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfu |
|---|---|---|---|---|---|---|
| **2942** | 3204 | B000084DVR | A1UGDJP1ZJWVPF | T. Moore "thoughtful reader" | 1 | 1 |

In [75]: `preprocessed_reviews[1500]`

Out[75]: 'wow far two two star reviews one obviously no idea ordering wants crispy cookies hey sorry reviews nobody good beyond reminding us look ordering chocolate oatmeal cookies not like combination not order type cookie find combo quite nice really oatmeal sort calms rich chocolate flavor gives cookie sort coconut type consistency let also remember tastes differ given opinion soft chewy cookies advertised not crispy cookies blurb would say crispy rather chewy happen like raw cookie dough however not see taste like raw cookie dough soft however confusion yes stick together soft cookies tend not individually wrapped would add cost oh yeah chocolate chip cookies tend somewhat sweet want something hard crisp suggest nabiso ginger snaps want cookie soft chewy tastes like combination chocolate oatmeal give try place second order'

## [3.2] Preprocess Summary

In [24]:
```
# printing some random reviews
sent_0 = final['Summary'].values[0]
print(sent_0)
print("="*50)

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

sent_1500 = final['Summary'].values[1500]
print(sent_1500)
```

```
print("="*50)

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

thirty bucks?
==================================================
Best sour cream & onion chip I've had
==================================================
Are We Reviewing Our Mistakes Or These Cookies?
==================================================
caribou
==================================================

In [25]:
```python
# remove urls from text python: https://stackoverflow.com/a/40823105/40
84039
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)
```

thirty bucks?

In [26]:
```python
# 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)
```

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

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

```
thirty bucks?
==================================================
Best sour cream & onion chip I've had
==================================================
Are We Reviewing Our Mistakes Or These Cookies?
==================================================
caribou
```

In [27]:
```python
# 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 [28]:
```python
sent_1500 = decontracted(sent_1500)
print(sent_1500)
```

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

```
Are We Reviewing Our Mistakes Or These Cookies?
==================================================
```

In [29]: 
```python
#remove words with numbers python: https://stackoverflow.com/a/1808237
0/4084039
sent_0 = re.sub("\S*\d\S*", "", sent_0).strip()
print(sent_0)
```

```
thirty bucks?
```

In [30]: 
```python
#remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)
print(sent_1500)
```

```
Are We Reviewing Our Mistakes Or These Cookies
```

In [31]: 
```python
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'no
t'
# <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', 'o
urs', 'ourselves', 'you', "you're", "you've",\
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselve
s', '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 [32]:
```python
# Combining all the above stundents
from tqdm import tqdm
preprocessed_Summary = []
# tqdm is for printing the status bar
for sentance in tqdm(final['Summary'].values):
    sentance = re.sub(r"http\S+", "", sentance)
    sentance = BeautifulSoup(sentance, 'lxml').get_text()
    sentance = decontracted(sentance)
    sentance = re.sub("\S*\d\S*", "", sentance).strip()
    sentance = re.sub('[^A-Za-z]+', ' ', sentance)
    # https://gist.github.com/sebleier/554280
    sentance = ' '.join(e.lower() for e in sentance.split() if e.lower
() not in stopwords)
    preprocessed_Summary.append(sentance.strip())
```
```
100%|███████████████████████████████| 4986/4986 [00:05<00:00, 91
6.80it/s]
```

In [33]:
```python
preprocessed_Summary[1500]
```

Out[33]: `'reviewing mistakes cookies'`

# [4] Featurization

## [4.1] BAG OF WORDS

In [76]:
```python
#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', 'abb
ott', '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 [77]:
```python
#bi-gram, tri-gram and n-gram

#removing stop words like "not" should be avoided before building n-gra
ms
# count_vect = CountVectorizer(ngram_range=(1,2))
# please do read the CountVectorizer documentation http://scikit-learn.
org/stable/modules/generated/sklearn.feature_extraction.text.CountVecto
rizer.html
# you can choose these numebrs min_df=10, max_features=5000, of your ch
oice
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_s
hape())
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 [78]:
```
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.ge
t_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', 'a
ble find', 'able get', 'absolute', 'absolutely', 'absolutely deliciou
s', '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

```python
In [79]:  # Train your own Word2Vec model using your own text corpus
          i=0
          list_of_sentance=[]
          for sentence in preprocessed_reviews:
              list_of_sentance.append(sentence.split())
```

```python
In [80]:  # 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 varible 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_sentance,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:
```

```python
    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_trai
n_w2v = True, to train your own w2v ")
```

```
[('excellent', 0.9955722093582153), ('especially', 0.9944643974304199),
('works', 0.9944315552711487), ('wonderful', 0.9943913221359253), ('gra
nola', 0.9943044781684875), ('also', 0.9940099716186523), ('general',
0.9938420057296753), ('quick', 0.9938058853149414), ('content', 0.99378
11493873596), ('watch', 0.9937753677368164)]
===================================================
[('oh', 0.9994736909866333), ('choice', 0.9994628429412842), ('looks',
0.9993342161178589), ('kernels', 0.9993302822113037), ('hands', 0.99931
53214454651), ('lover', 0.9993065595626831), ('opinion', 0.999306499958
0383), ('berry', 0.9992810487747192), ('device', 0.9992711544036865),
('stash', 0.9992539286613464)]
```

In [81]:
```python
w2v_words = list(w2v_model.wv.vocab)
print("number of words that occured minimum 5 times ",len(w2v_words))
print("sample words ", w2v_words[0:50])
```

```
number of words that occured minimum 5 times  3817
sample words  ['product', 'available', 'course', 'total', 'pretty', 'st
inky', 'right', 'nearby', 'used', 'ca', 'not', 'beat', 'great', 'receiv
ed', 'shipment', 'could', 'hardly', 'wait', 'try', 'love', 'call', 'ins
tead', 'removed', 'easily', 'daughter', 'designed', 'printed', 'use',
'car', 'windows', 'beautifully', 'shop', 'program', 'going', 'lot', 'fu
n', 'everywhere', 'like', 'tv', 'computer', 'really', 'good', 'idea',
'final', 'outstanding', 'window', 'everybody', 'asks', 'bought', 'mad
e']
```

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

**[4.4.1.1] Avg W2v**

In [82]:
```python
# 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_sentance): # 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 [00:16<00:00, 29
7.43it/s]
```

```
4986
50
```

**[4.4.1.2] TFIDF weighted W2v**

In [83]:
```python
# 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 v
alue
dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
```

In [84]:
```python
# 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 ce
ll_val = tfidf

tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is st
ored in this list
row=0;
for sent in tqdm(list_of_sentance): # 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/r
eview
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
#             tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    tfidf_sent_vectors.append(sent_vec)
    row += 1
```

```
100%|████████████████████████████| 4986/4986 [01:28<00:00, 5
6.11it/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)

## [5.1] Applying TNSE on Text BOW vectors

```
In [86]: # please write all the code with proper documentation, and proper title
         s 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
         final_counts.shape
```

```
Out[86]: (4986, 12997)
```

```
In [87]: # Change sparse matrix to dense matrix
         final_counts = final_counts.todense()
```

```
In [88]: import warnings
         warnings.filterwarnings('ignore')
         # Data-preprocessing: Standardizing the data

         from sklearn.preprocessing import StandardScaler
         standardized_data = StandardScaler().fit_transform(final_counts)
         print(standardized_data.shape)
```

```
(4986, 12997)
```

```
In [89]: final['Score'].value_counts()
```

```
Out[89]: 1    4178
         0     808
         Name: Score, dtype: int64
```

```
In [91]: # TSNE on BOW vectors
```

```python
from sklearn.manifold import TSNE


model = TSNE(n_components=2, random_state=0)
# configuring the parameteres
# the number of components = 2
# default perplexity = 30
# default learning rate = 200
# default Maximum number of iterations for the optimization = 1000

tsne_data = model.fit_transform(standardized_data)


# creating a new data frame which help us in ploting the result data
tsne_data = np.vstack((tsne_data.T, final['Score'])).T
tsne_df = pd.DataFrame(data=tsne_data, columns=("Dim_1", "Dim_2", "labe
l"))

# Ploting the result of tsne
sns.FacetGrid(tsne_df, hue="label", size=8).map(plt.scatter, 'Dim_1',
'Dim_2').add_legend()
plt.title('t-SNE of Bag of Words(BoW) with perplexity = 30 and n_iter =
 1000',size=20)
plt.show()
```
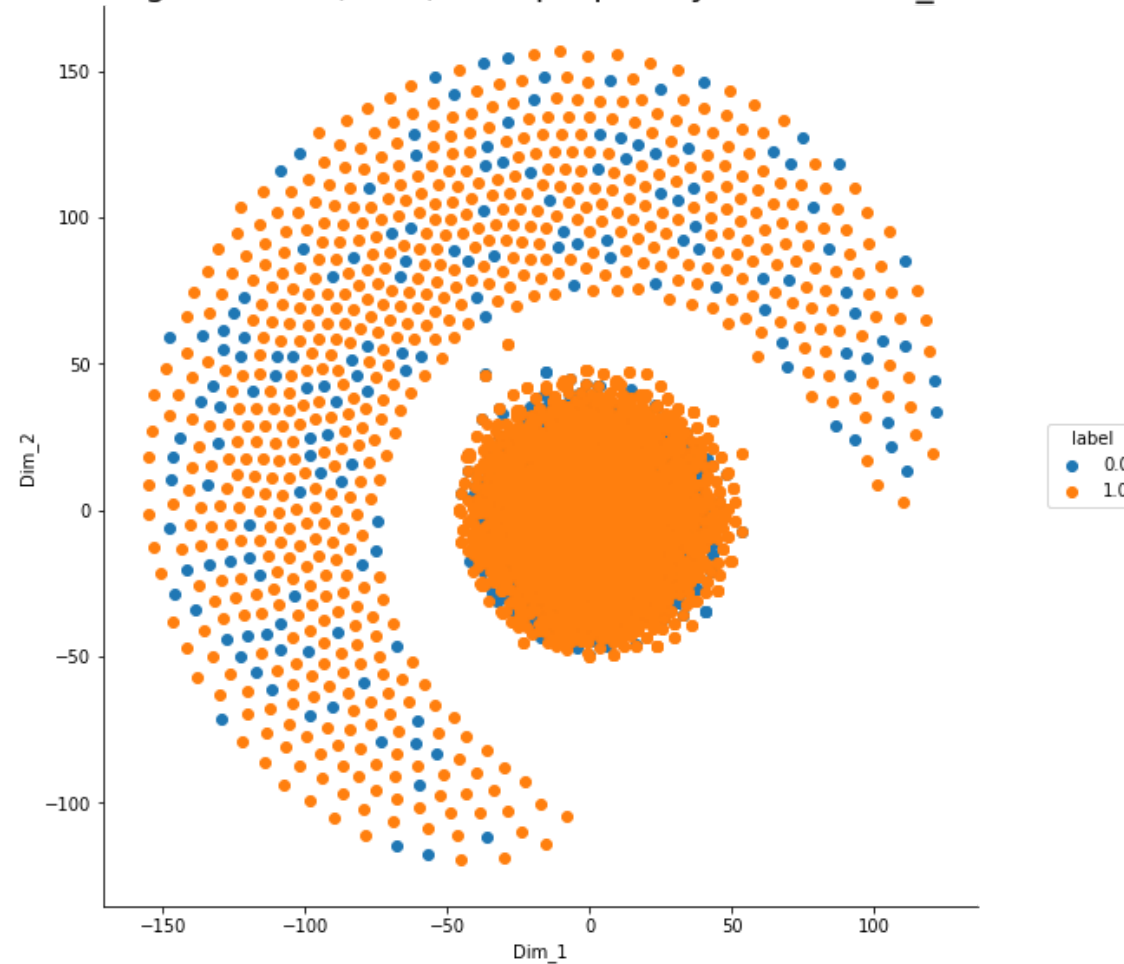
t-SNE of Bag of Words(BoW) with perplexity = 30 and n_iter = 1000

In [92]:
```python
# TSNE

from sklearn.manifold import TSNE

model = TSNE(n_components=2, random_state=0, perplexity=50,  n_iter=200
0)
```

```python
# configuring the parameteres
# the number of components = 2
# default perplexity = 30
# default learning rate = 200
# default Maximum number of iterations for the optimization = 1000

tsne_data = model.fit_transform(standardized_data)


# creating a new data frame which help us in ploting the result data
tsne_data = np.vstack((tsne_data.T, final['Score'])).T
tsne_df = pd.DataFrame(data=tsne_data, columns=("Dim_1", "Dim_2", "label"))

# Ploting the result of tsne
sns.FacetGrid(tsne_df, hue="label", size=8).map(plt.scatter, 'Dim_1', 'Dim_2').add_legend()
plt.title('t-SNE of Bag of Words(BoW) with perplexity = 50 and n_iter = 2000',size=20)
plt.show()
```

t-SNE of Bag of Words(BoW) with perplexity = 30 and n_iter = 1000

## observation

- 1.with increasing number of iterations and perplexity overlapping of points are dense on one another i.e +ve points are overlapped by -ve points.
- 2.so,it is difficult of a line to classify the points

## [5.1] Applying TNSE on Text TFIDF vectors

```
In [93]:  # please write all the code with proper documentation, and proper title
          s 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
          features = tf_idf_vect.get_feature_names()
          print("some sample features(unique words in the corpus)",features[1000:
          1010])
```

```
some sample features(unique words in the corpus) ['food', 'food allergi
es', 'food dog', 'food good', 'food items', 'food like', 'food no', 'fo
od not', 'food one', 'food store']
```

```
In [94]:  # Change sparse matrix to dense matrix
          final_tf_idf  = final_tf_idf.todense()
```

```
In [95]:  final_tf_idf.shape
```

```
Out[95]:  (4986, 3144)
```

```
In [96]:  import warnings
          warnings.filterwarnings('ignore')
          # Data-preprocessing: Standardizing the data

          from sklearn.preprocessing import StandardScaler
          standardized_data = StandardScaler().fit_transform(final_tf_idf)
          print(standardized_data.shape)
```
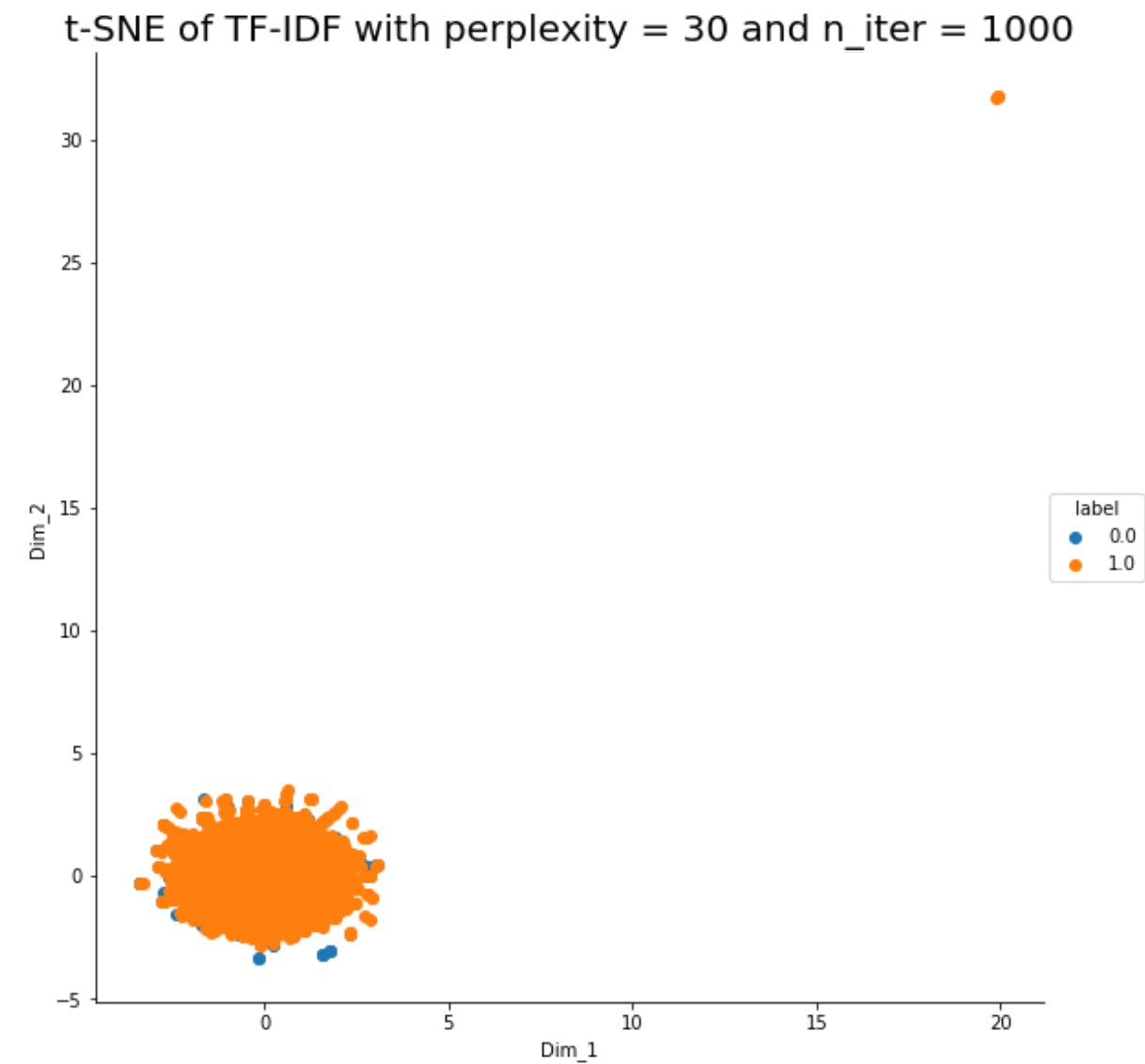
```
(4986, 3144)
```

```
In [97]:  # TSNE on text TFIDF vectors
```

```python
from sklearn.manifold import TSNE

# t-SNE with perplexity = 50 and n_iter = 3000
model = TSNE(n_components=2, random_state=0, perplexity=50,  n_iter=300
0)

tsne_data = model.fit_transform(standardized_data)


# creating a new data frame which help us in ploting the result data
tsne_data = np.vstack((tsne_data.T, final['Score'])).T
tsne_df = pd.DataFrame(data=tsne_data, columns=("Dim_1", "Dim_2", "labe
l"))

# Ploting the result of tsne
sns.FacetGrid(tsne_df, hue="label", size=8).map(plt.scatter, 'Dim_1',
'Dim_2').add_legend()
plt.title('t-SNE of TF-IDF with perplexity = 50 and n_iter = 3000',size
=20)
plt.show()
```

t-SNE of TF-IDF with perplexity = 50 and n_iter = 3000

In [98]:
```python
# TSNE on text TFIDF vectors

from sklearn.manifold import TSNE
```

```python
# t-SNE with perplexity = 30 and n_iter = 1000
model = TSNE(n_components=2, random_state=0, perplexity=30,  n_iter=100
0)

tsne_data = model.fit_transform(standardized_data)


# creating a new data frame which help us in ploting the result data
tsne_data = np.vstack((tsne_data.T, final['Score'])).T
tsne_df = pd.DataFrame(data=tsne_data, columns=("Dim_1", "Dim_2", "labe
l"))

# Ploting the result of tsne
sns.FacetGrid(tsne_df, hue="label", size=8).map(plt.scatter, 'Dim_1',
'Dim_2').add_legend()
plt.title('t-SNE of TF-IDF with perplexity = 30 and n_iter = 1000',size
=20)
plt.show()
```

t-SNE of TF-IDF with perplexity = 30 and n_iter = 1000

**observation**

- 1.observing above plots we conclude that as the perplexity and number of iterations increases the overlapping of both the classes decreases and also the density of classes around the plot tend to decrease and then Increased in later increase of perplexity and iterations.

## [5.3] Applying TNSE on Text Avg W2V vectors

In [99]:
```
# please write all the code with proper documentation, and proper title
s 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
import warnings
warnings.filterwarnings('ignore')
# Data-preprocessing: Standardizing the data

from sklearn.preprocessing import StandardScaler
standardized_data = StandardScaler().fit_transform(sent_vectors)
print(standardized_data.shape)
```

```
(4986, 50)
```

In [100]:
```
# TSNE on AVG w2v

from sklearn.manifold import TSNE

# t-SNE with perplexity = 30 and n_iter = 500
model = TSNE(n_components=2, random_state=0, perplexity=30,  n_iter=500
)

tsne_data = model.fit_transform(standardized_data)


# creating a new data frame which help us in ploting the result data
```
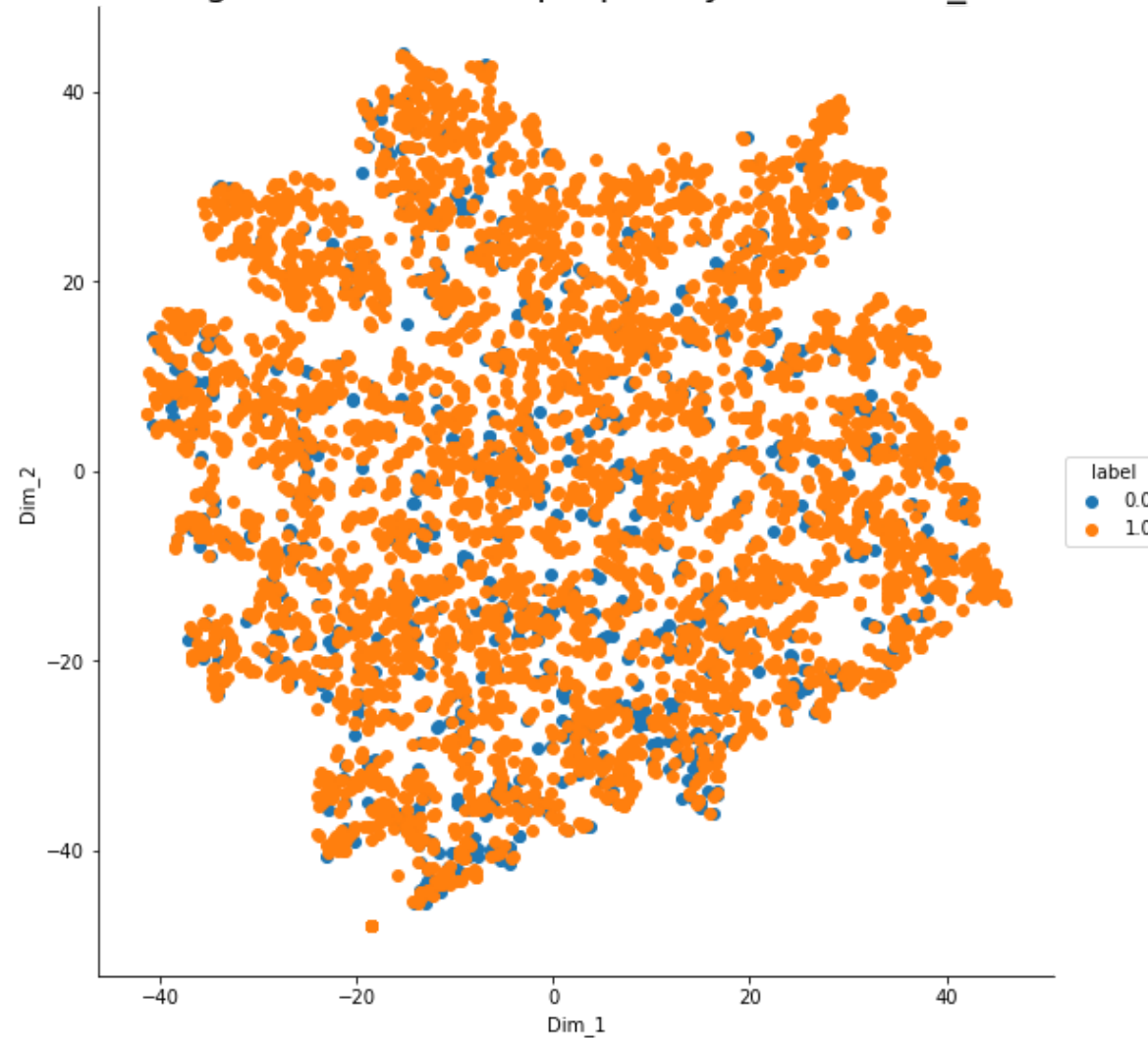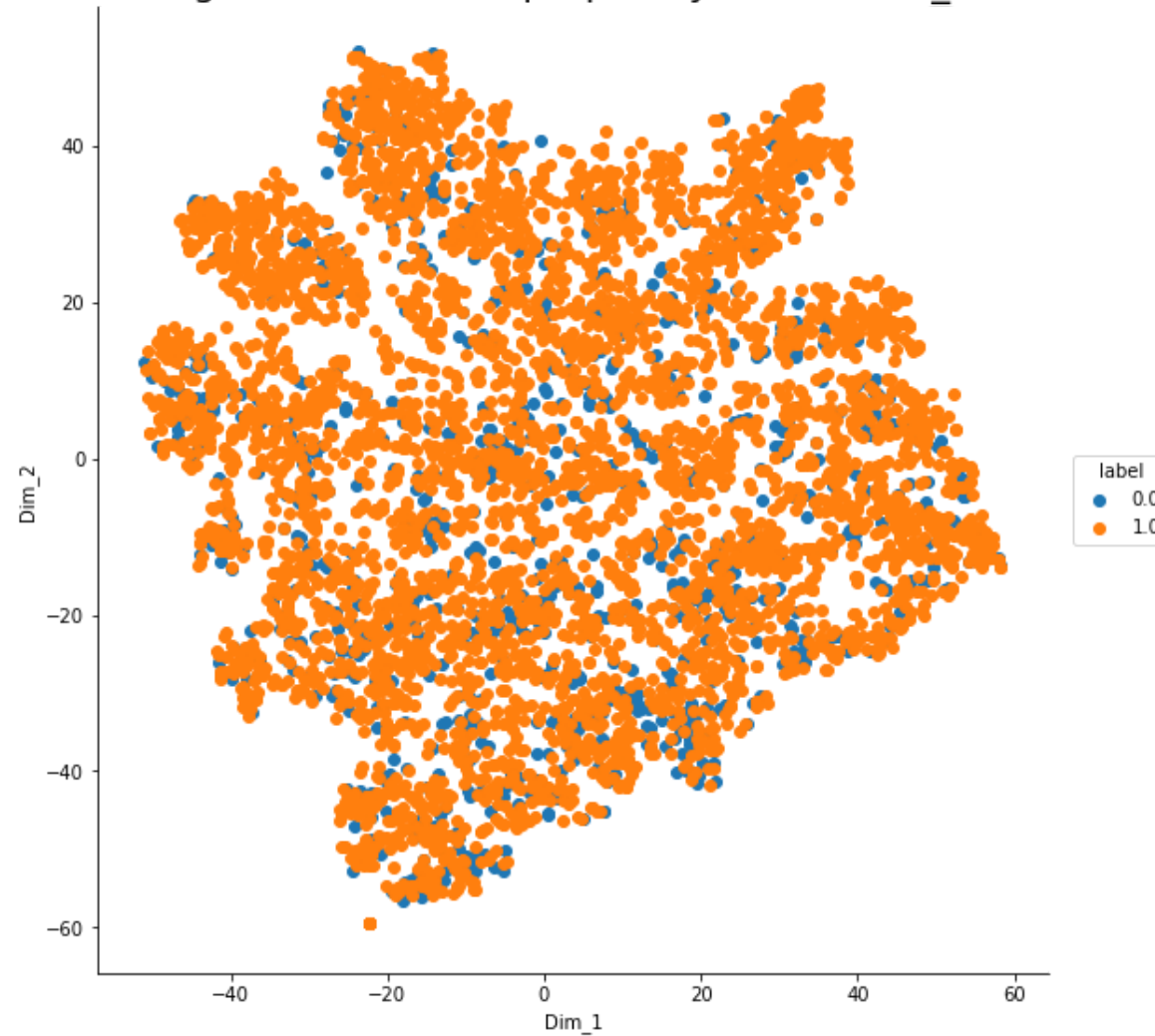
```python
tsne_data = np.vstack((tsne_data.T, final['Score'])).T
tsne_df = pd.DataFrame(data=tsne_data, columns=("Dim_1", "Dim_2", "labe
l"))

# Ploting the result of tsne
sns.FacetGrid(tsne_df, hue="label", size=8).map(plt.scatter, 'Dim_1',
'Dim_2').add_legend()
plt.title('t-SNE of Avg Word2Vec with perplexity = 30 and n_iter = 500'
,size=20)
plt.show()
```

## t-SNE of Avg Word2Vec with perplexity = 30 and n_iter = 500



In [101]:
```python
# TSNE on avg w2v

from sklearn.manifold import TSNE

# t-SNE with perplexity = 50 and n_iter = 1000
```

```python
model = TSNE(n_components=2, random_state=0, perplexity=50,  n_iter=100
0)

tsne_data = model.fit_transform(standardized_data)


# creating a new data frame which help us in ploting the result data
tsne_data = np.vstack((tsne_data.T, final['Score'])).T
tsne_df = pd.DataFrame(data=tsne_data, columns=("Dim_1", "Dim_2", "labe
l"))

# Ploting the result of tsne
sns.FacetGrid(tsne_df, hue="label", size=8).map(plt.scatter, 'Dim_1',
'Dim_2').add_legend()
plt.title('t-SNE of Avg Word2Vec with perplexity = 50 and n_iter = 100
0',size=20)
plt.show()
```

t-SNE of Avg Word2Vec with perplexity = 50 and n_iter = 1000

## observation

- observing above plots we conclude that as the perplexity and number of iterations increases the area of covered by the classes on the plot decreased.and also difficult to

clasify .

## [5.4] Applying TNSE on Text TFIDF weighted W2V vectors

In [102]:
```python
# please write all the code with proper documentation, and proper title
s 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
import warnings
warnings.filterwarnings('ignore')
# Data-preprocessing: Standardizing the data

from sklearn.preprocessing import StandardScaler
standardized_data = StandardScaler().fit_transform(tfidf_sent_vectors)
print(standardized_data.shape)
```

(4986, 50)

In [103]:
```python
# TSNE with TFIDF weighted W2V

from sklearn.manifold import TSNE

# t-SNE with perplexity = 30 and n_iter = 500
model = TSNE(n_components=2, random_state=0, perplexity=30,  n_iter=500
)

tsne_data = model.fit_transform(standardized_data)


# creating a new data frame which help us in ploting the result data
tsne_data = np.vstack((tsne_data.T, final['Score'])).T
tsne_df = pd.DataFrame(data=tsne_data, columns=("Dim_1", "Dim_2", "labe
```
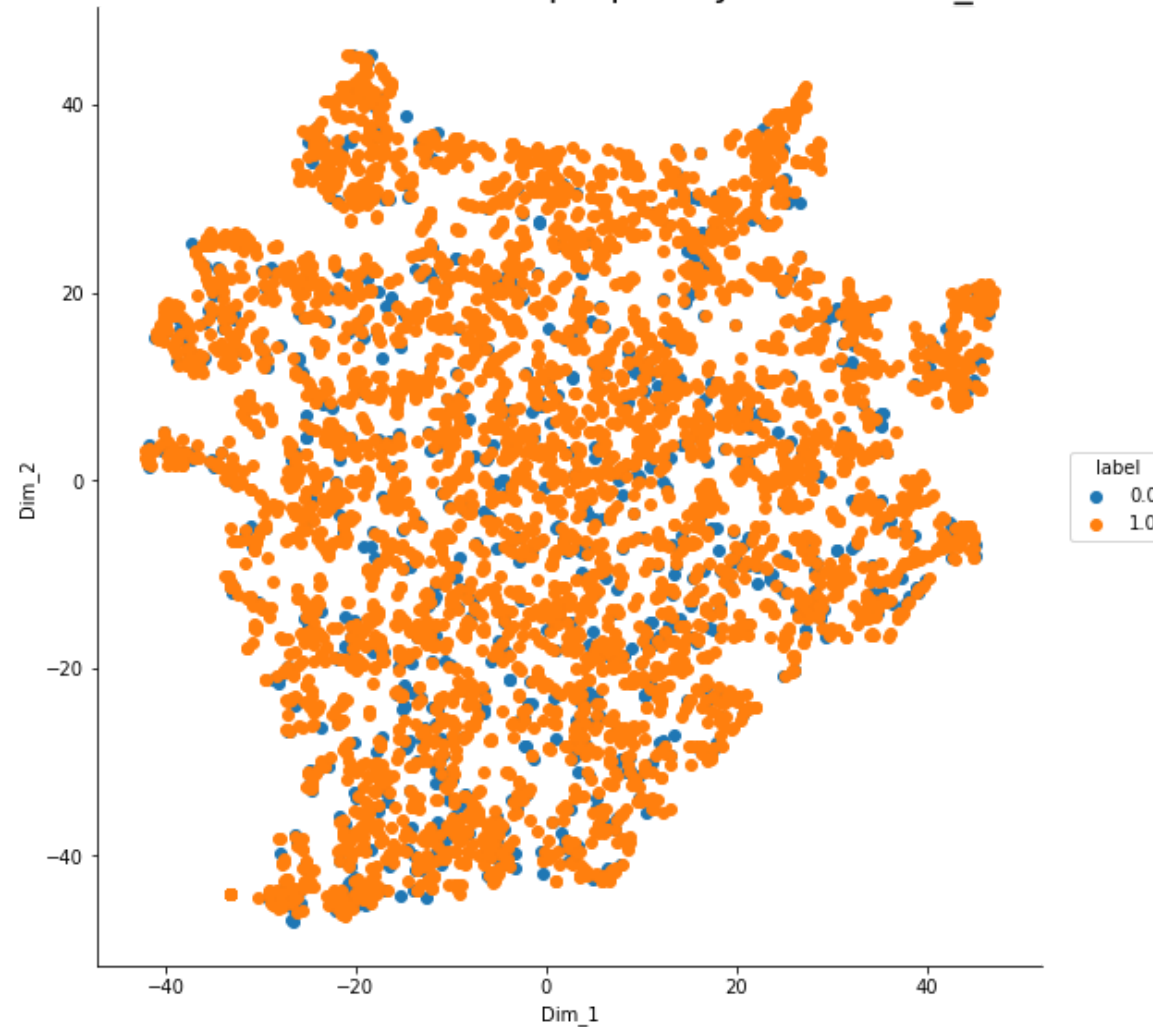
```python
l"))

# Ploting the result of tsne
sns.FacetGrid(tsne_df, hue="label", size=8).map(plt.scatter, 'Dim_1',
'Dim_2').add_legend()
plt.title('t-SNE of TFIDF Word2Vec with perplexity = 30 and n_iter = 50
0',size=20)
plt.show()
```

t-SNE of TFIDF Word2Vec with perplexity = 30 and n_iter = 500

t-SNE of TFIDF Word2vec with perplexity = 30 and n_iter = 500

In [105]:
```python
# TSNE with TFIDF weighted W2V

from sklearn.manifold import TSNE

# t-SNE with perplexity = 30 and n_iter = 500
model = TSNE(n_components=2, random_state=0, perplexity=50,  n_iter=100
0)
```
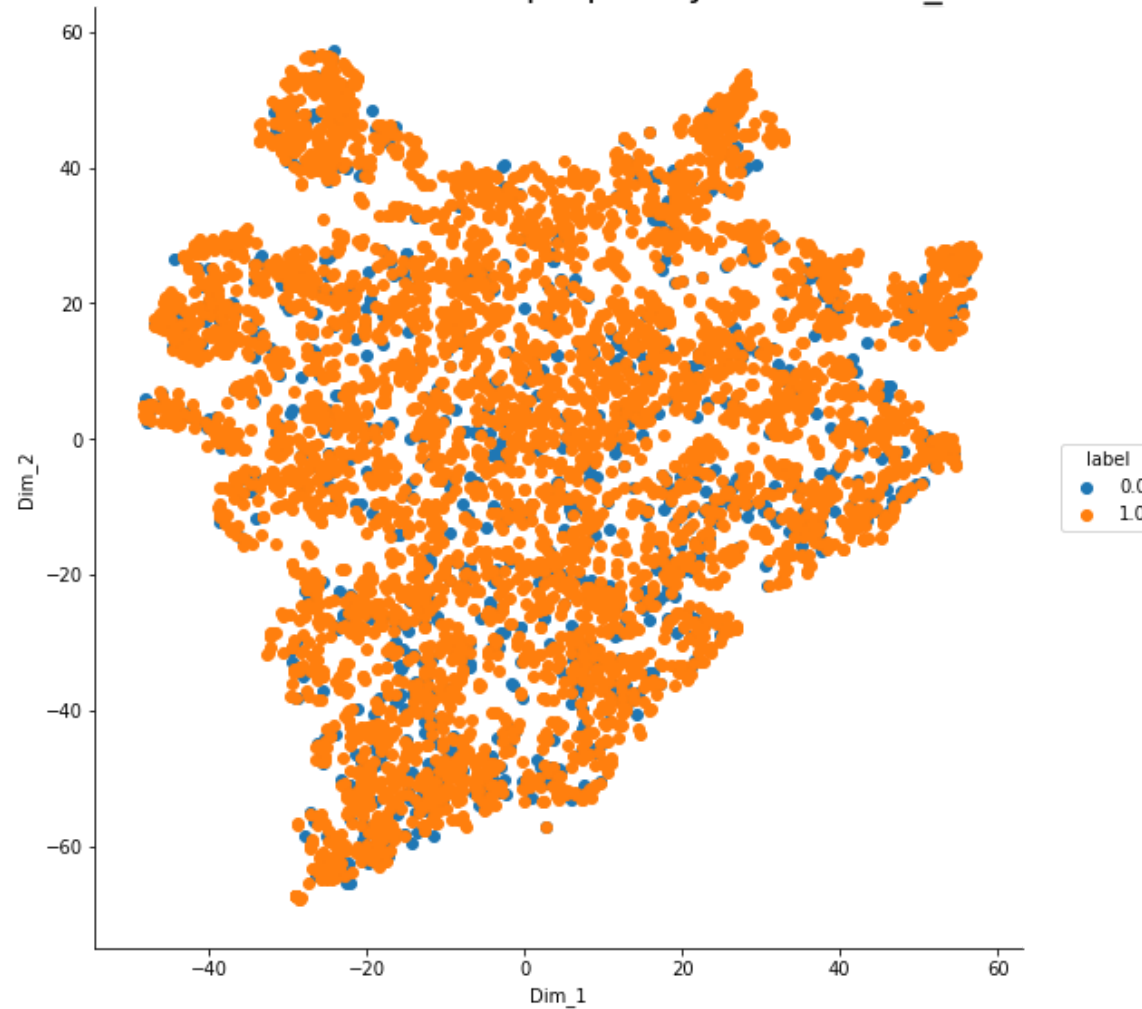
```python
tsne_data = model.fit_transform(standardized_data)


# creating a new data frame which help us in ploting the result data
tsne_data = np.vstack((tsne_data.T, final['Score'])).T
tsne_df = pd.DataFrame(data=tsne_data, columns=("Dim_1", "Dim_2", "labe
l"))

# Ploting the result of tsne
sns.FacetGrid(tsne_df, hue="label", size=8).map(plt.scatter, 'Dim_1',
'Dim_2').add_legend()
plt.title('t-SNE of TFIDF Word2Vec with perplexity = 50 and n_iter = 10
00',size=20)
plt.show()
```

t-SNE of TFIDF Word2Vec with perplexity = 50 and n_iter = 1000

## observation

- observing above plots we conclude that as the perplexity and number of iterations increases the overlapping of both the classes also increases

# [6] Conclusions

- BOW->with increasing number of iterations and perplexity overlapping of points are dense on one another i.e +ve points are overlapped by -ve points.so,it is difficult of a line to classify the points
- TF-IDF->observing above plots we conclude that as the perplexity and number of iterations increases the overlapping of both the classes decreases and also the density of classes around the plot tend to decrease and then Increased in later increase of perplexity and iterations.
- AVG W2V->observing above plots we conclude that as the perplexity and number of iterations increases the area of covered by the classes on the plot decreased.and also difficult to clasify .
- TFIDF WEG W2V->observing above plots we conclude that as the perplexity and number of iterations increases the overlapping of both the classes alsp increases