

# TSNE-amazon-fine-food-reviews

July 30, 2018

## 1 [7] Amazon Fine Food Reviews Analysis

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

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.

### 1.1 [7.1] Loading the data

The dataset is available in two forms 1. .csv file 2. SQLite Database

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

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

```
In [1]: %matplotlib inline
```

```
import 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
```

```
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
```

```
filtered_data = pd.read_sql_query("""SELECT * FROM Reviews WHERE Score != 3""", 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 'negative'
```

```
    return 'positive'
```

```
#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
```

```
In [3]: filtered_data.shape #looking at the number of attributes and size of the data
```

```
filtered_data.head()
```

```
Out[3]:
```

	Id	ProductId	UserId	ProfileName	\
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	
2	3	B000LQOCHO	ABXLMWJIXXAIN	Natalia Corres	"Natalia Corres"
3	4	B000UA0QIQ	A395BORC6FGVXV	Karl	
4	5	B006K2ZZ7K	A1UQRSCLF8GW1T	Michael D. Bigham	"M. Wassir"

	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	\
0	1	1	positive	1303862400	

1	0	0	negative	1346976000
2	1	1	positive	1219017600
3	3	3	negative	1307923200
4	0	0	positive	1350777600

	Summary	Text
0	Good Quality Dog Food	I have bought several of the Vitality canned d...
1	Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...
2	"Delight" says it all	This is a confection that has been around a fe...
3	Cough Medicine	If you are looking for the secret ingredient i...
4	Great taffy	Great taffy at a great price. There was a wid...

## 2 Exploratory Data Analysis

### 2.1 [7.1.2] 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 [5]: display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID
""", con)
display
```

```
Out[5]:
```

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator \
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

	HelpfulnessDenominator	Score	Time \
0	2	5	1199577600
1	2	5	1199577600
2	2	5	1199577600
3	2	5	1199577600
4	2	5	1199577600

	Summary \
0	LOACKER QUADRATINI VANILLA WAFERS
1	LOACKER QUADRATINI VANILLA WAFERS
2	LOACKER QUADRATINI VANILLA WAFERS
3	LOACKER QUADRATINI VANILLA WAFERS
4	LOACKER QUADRATINI VANILLA WAFERS

	Text
0	DELICIOUS WAFERS. I FIND THAT EUROPEAN WAFERS ...
1	DELICIOUS WAFERS. I FIND THAT EUROPEAN WAFERS ...
2	DELICIOUS WAFERS. I FIND THAT EUROPEAN WAFERS ...
3	DELICIOUS WAFERS. I FIND THAT EUROPEAN WAFERS ...
4	DELICIOUS WAFERS. I FIND THAT EUROPEAN WAFERS ...

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 [6]: #Sorting data according to ProductId in ascending order
```

```
sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True, inplace=False)
```

```
In [7]: #Deduplication of entries
```

```
final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"}, keep='first')
final.shape
```

```
Out[7]: (364173, 10)
```

```
In [8]: #Checking to see how much % of data still remains
```

```
(final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

```
Out[8]: 69.25890143662969
```

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

```

Out[9]:      Id  ProductId      UserId      ProfileName \
0  64422  B000MIDROQ  A161DK06JJMCYF  J. E. Stephens "Jeanne"
1  44737  B001EQ55RW  A2VOI904FH7ABY                      Ram

      HelpfulnessNumerator  HelpfulnessDenominator  Score      Time \
0                        3                        1      5  1224892800
1                        3                        2      4  1212883200

                        Summary \
0          Bought This for My Son at College
1  Pure cocoa taste with crunchy almonds inside

                        Text
0  My son loves spaghetti so I didn't hesitate or...
1  It was almost a 'love at first bite' - the per...

```

```
In [10]: final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [11]: #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()

```

```
(364171, 10)
```

```

Out[11]: positive    307061
        negative     57110
        Name: Score, dtype: int64

```

## 2.2 7.2.3 Text Preprocessing: Stemming, stop-word removal and Lemmatization.

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 [12]: *# find sentences containing HTML tags*

```
import re
i=0;
for sent in final['Text'].values:
    if (len(re.findall('<.*?>', sent))):
        print(i)
        print(sent)
        break;
    i += 1;
```

6

I set aside at least an hour each day to read to my son (3 y/o). At this point, I consider mys

In [13]: *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
```

```
stop = set(stopwords.words('english')) #set of stopwords
```

```
sno = nltk.stem.SnowballStemmer('english') #initialising the snowball stemmer
```

```
def cleanhtml(sentence): #function to clean the word of any html-tags
```

```
    cleanr = re.compile('<.*?>')
    cleantext = re.sub(cleanr, ' ', sentence)
    return cleantext
```

```
def cleanpunc(sentence): #function to clean the word of any punctuation or special ch
```

```
    cleaned = re.sub(r'[?|!|\\'|"|#]',r'',sentence)
    cleaned = re.sub(r'[,|,|)|(|\\|/]',r'',cleaned)
    return cleaned
```

```
print(stop)
```

```
print('*****')
```

```
print(sno.stem('tasty'))
```

{'in', 'above', 'hadn', 'myself', 'her', 'herself', 'yourselves', 'own', 'at', 'wasn', "you're"

\*\*\*\*\*

tasti

In [14]: *#Code for implementing step-by-step the checks mentioned in the pre-processing phase*

*# this code takes a while to run as it needs to run on 500k sentences.*

```
i=0
```

```
str1=' '
```

```
final_string=[]
```

```
all_positive_words=[] # store words from +ve reviews here
```

```

all_negative_words=[] # store words from -ve reviews here.
s=''
for sent in final['Text'].values:
    filtered_sentence=[]
    #print(sent);
    sent=cleanhtml(sent) # remove HTML tags
    for w in sent.split():
        for cleaned_words in cleanpunc(w).split():
            if((cleaned_words.isalpha()) & (len(cleaned_words)>2)):
                if(cleaned_words.lower() not in stop):
                    s=(sno.stem(cleaned_words.lower())).encode('utf8')
                    filtered_sentence.append(s)
                    if (final['Score'].values[i] == 'positive':
                        all_positive_words.append(s) #list of all words used to descr
                    if(final['Score'].values[i] == 'negative':
                        all_negative_words.append(s) #list of all words used to descr
                else:
                    continue
            else:
                continue
    #print(filtered_sentence)
    str1 = b" ".join(filtered_sentence) #final string of cleaned words
    #print("*****")

    final_string.append(str1)
    i+=1

```

In [15]: final['CleanedText']=final\_string #adding a column of CleanedText which displays the

In [16]: final.head(3) #below the processed review can be seen in the CleanedText Column

```

# store final table into an SQLite table for future.
conn = sqlite3.connect('final.sqlite')
c=conn.cursor()
conn.text_factory = str
final.to_sql('Reviews', conn, flavor=None, schema=None, if_exists='replace', index=True)

```

### 3 [7.2.2] Bag of Words (BoW)

In [17]: # To get 2k +ve and 2k -ve reviews randomly as my system has only 4gb RAM.  
# It gives MemoryError when I take more than 4k review.

```

data_pos = final[final["Score"] == "positive"].sample(n = 2000)
data_neg = final[final["Score"] == "negative"].sample(n = 2000)
final_4000 = pd.concat([data_pos, data_neg])

```

In [18]: #final\_4000

```

In [19]: score_4000 = final_4000["Score"]

In [20]: score_4000.shape

Out[20]: (4000,)

In [21]: final_4000.shape

Out[21]: (4000, 11)

In [22]: #BoW
          count_vect = CountVectorizer() #in scikit-learn
          final_counts = count_vect.fit_transform(final_4000['CleanedText'].values)

In [23]: #final_counts

In [24]: type(final_counts)

Out[24]: scipy.sparse.csr.csr_matrix

In [25]: final_counts.get_shape()

Out[25]: (4000, 8714)

```

### 3.1 [7.2.4] Bi-Grams and n-Grams.

#### Motivation

Now that we have our list of words describing positive and negative reviews lets analyse them. We begin analysis by getting the frequency distribution of the words as shown below

```

In [26]: freq_dist_positive=nlTK.FreqDist(all_positive_words)
          freq_dist_negative=nlTK.FreqDist(all_negative_words)
          print("Most Common Positive Words : ",freq_dist_positive.most_common(20))
          print("Most Common Negative Words : ",freq_dist_negative.most_common(20))

```

```

Most Common Positive Words : [(b'like', 139429), (b'tast', 129047), (b'good', 112766), (b'fla
Most Common Negative Words : [(b'tast', 34585), (b'like', 32330), (b'product', 28218), (b'one

```

Observation:- From the above it can be seen that the most common positive and the negative words overlap for eg. 'like' could be used as 'not like' etc. So, it is a good idea to consider pairs of consequent words (bi-grams) or q sequence of n consecutive words (n-grams)

```

In [27]: #bi-gram, tri-gram and n-gram

          #removing stop words like "not" should be avoided before building n-grams
          count_vect = CountVectorizer(ngram_range=(1,2) ) #in scikit-learn
          final_bigram_counts = count_vect.fit_transform(final_4000['CleanedText'].values)

In [28]: final_bigram_counts.get_shape()

```



Out[28]: (4000, 115573)

```
In [29]: from sklearn.preprocessing import StandardScaler
```

```
std_data = StandardScaler(with_mean = False).fit_transform(final_bigram_counts)
std_data.shape
```

```
C:\Users\premvardhan\Anaconda3\lib\site-packages\sklearn\utils\validation.py:429: DataConversionWarning:
  warnings.warn(msg, _DataConversionWarning)
```

Out[29]: (4000, 115573)

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

Out[30]: scipy.sparse.csr.csr\_matrix

```
In [31]: # convert sparse to dense as tsne takes dense vector
std_data = std_data.todense()
```

```
In [32]: type(std_data)
```

Out[32]: numpy.matrixlib.defmatrix.matrix

```
In [33]: from sklearn.manifold import TSNE
model = TSNE(n_components=2, random_state=0, perplexity = 30, n_iter = 5000)
# 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(std_data)
```

```
# creating a new data frame which help us in plotting the result data
```

```
tsne_data = np.vstack((tsne_data.T, score_4000)).T
```

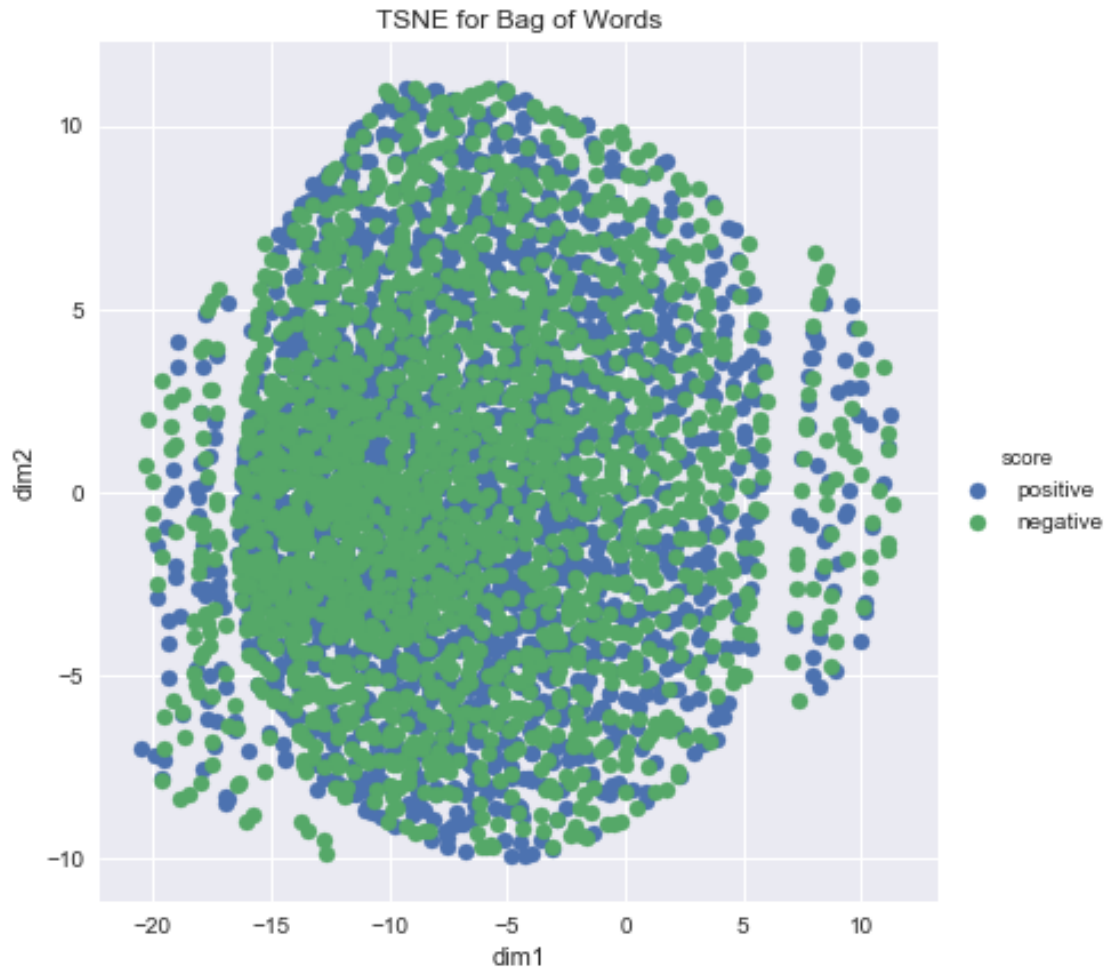
```
tsne_df = pd.DataFrame(data=tsne_data, columns=("dim1", "dim2", "score"))
```

```
# Ploting the result of tsne
```

```
sns.FacetGrid(tsne_df, hue="score", size=6).map(plt.scatter, 'dim1', 'dim2').add_legend
```

```
plt.title("TSNE for Bag of Words")
```

```
plt.show()
```



**Observation:-** Here, we are unable to simply draw a hyperplane and separate +ve and -ve reviews because they overlap each other. But we will have some alternative way to separate reviews.

## 4 [7.2.5] TF-IDF

```
In [34]: # Tf-Idf
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2))
final_tf_idf = tf_idf_vect.fit_transform(final_4000['CleanedText'].values)

In [35]: #final_4000

In [36]: # Standardization
from sklearn.preprocessing import StandardScaler
std = StandardScaler(with_mean = False)
std_data = std.fit_transform(final_tf_idf)

In [37]: # Converting sparse matrix to dense because tnse takes dense vector
std_data = std_data.todense()
```

```
In [38]: std_data.shape
```

```
Out[38]: (4000, 115573)
```

```
In [39]: # tsne
```

```
from sklearn.manifold import TSNE
```

```
model = TSNE(n_components = 2, perplexity = 50)
```

```
tsne_data = model.fit_transform(std_data)
```

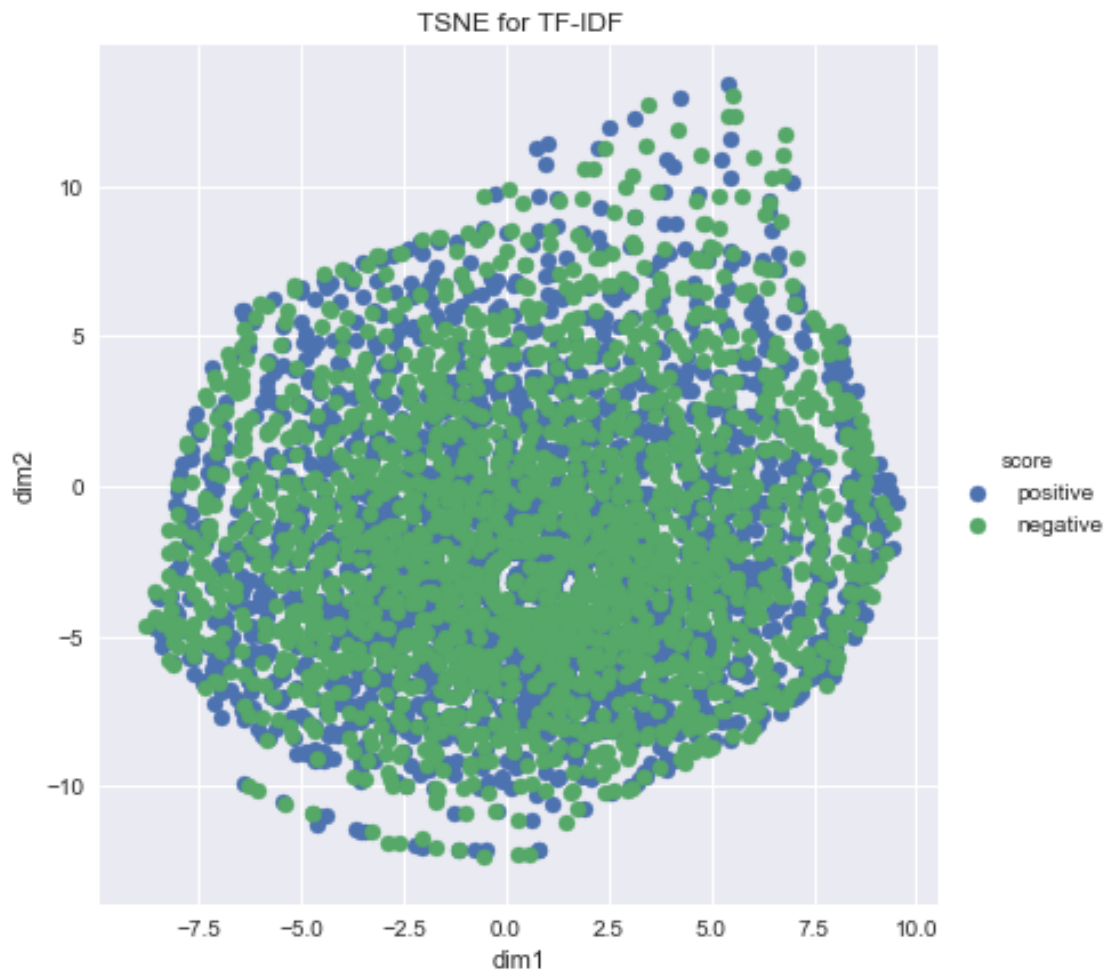
```
tsne_data = np.vstack((tsne_data.T, score_4000)).T
```

```
tsne_df = pd.DataFrame(data = tsne_data, columns = ("dim1", "dim2", "score"))
```

```
sns.FacetGrid(tsne_df, hue = "score", size = 6).map(plt.scatter, "dim1", "dim2").add_
```

```
plt.title("TSNE for TF-IDF")
```

```
plt.show()
```



```
In [40]: features = tf_idf_vect.get_feature_names()  
len(features)
```

```
Out[40]: 115573
```

```
In [41]: # convert a row in sparse matrix to a numpy array
         print(final_tf_idf[3,:].toarray()[0])
```

```
[ 0.  0.  0. ...,  0.  0.  0.]
```

```
In [42]: # source: https://buhrmann.github.io/tfidf-analysis.html
         def top_tfidf_feats(row, features, top_n=25):
             ''' Get top n tfidf values in row and return them with their corresponding features
             topn_ids = np.argsort(row)[:top_n]
             top_feats = [(features[i], row[i]) for i in topn_ids]
             df = pd.DataFrame(top_feats)
             df.columns = ['feature', 'tfidf']
             return df
```

```
         top_tfidf = top_tfidf_feats(final_tf_idf[1,:].toarray()[0], features, 25)
```

```
In [43]: #top_tfidf
```

**Observations:-** As this representation also looks like bow and massively overlapped +ve and -ve review.

## 5 [7.2.6] Word2Vec

```
In [ ]: # Using Google News Word2Vectors
         from gensim.models import Word2Vec
         from gensim.models import KeyedVectors
         import pickle
```

```
         # 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 which contains a dict ,
         # and it contains all our corpus words as keys and model[word] as values
         # To use this code-snippet, download "GoogleNews-vectors-negative300.bin"
         # from https://drive.google.com/file/d/0B7XkCwpI5KDYNlNUTTlSS21pQmM/edit
         # it's 1.9GB in size.
```

```
         #model = KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin', binary=True)
```

```
In [45]: #model.wv['computer']
```

```
In [46]: #model.wv.similarity('woman', 'man')
```

```
In [47]: #model.wv.most_similar('woman')
```

```
In [48]: #model.wv.most_similar('tasti') # "tasti" is the stemmed word for tasty, tastful
```

```
In [49]: #model.wv.most_similar('tasty')
```

```
In [50]: #model.wv.similarity('tasty', 'tast')
```

```
In [51]: # Train your own Word2Vec model using your own text corpus
```

```
import gensim
list_of_sent = []
for sent in final_4000['Text'].values:
    filtered_sentence = []
    sent=cleanhtml(sent)
    for w in sent.split():
        for cleaned_words in cleanpunc(w).split():
            if(cleaned_words.isalpha()):
                filtered_sentence.append(cleaned_words.lower())
            else:
                continue
    list_of_sent.append(filtered_sentence)
```

```
In [52]: print(final_4000['Text'].values[0])
print("*****")
print(list_of_sent[0])
```

These peanuts have just enough salt to keep them from tasting bland, but not enough to cause y  
\*\*\*\*\*  
['these', 'peanuts', 'have', 'just', 'enough', 'salt', 'to', 'keep', 'them', 'from', 'tasting']

```
In [53]: w2v_model=gensim.models.Word2Vec(list_of_sent,min_count=5,size=50, workers=4)
```

```
In [54]: w2v = w2v_model[w2v_model.wv.vocab]
```

C:\Users\premvardhan\Anaconda3\lib\site-packages\ipykernel\_launcher.py:1: DeprecationWarning: (  
 """Entry point for launching an IPython kernel.

```
In [55]: w2v.shape
```

```
Out[55]: (3854, 50)
```

```
In [56]: #w2v_2000 = w2v[0:2000:]
#w2v_2000.shape
```

```
In [57]: words = list(w2v_model.wv.vocab)
print(len(words))
```

3854

```

In [58]: w2v_model.wv.most_similar('tasty')

Out[58]: [('delicious', 0.9884324669837952),
          ('pretty', 0.9779555797576904),
          ('overly', 0.9710423350334167),
          ('healthy', 0.9650558233261108),
          ('salty', 0.9634706974029541),
          ('awfully', 0.9603949189186096),
          ('spicy', 0.9586955904960632),
          ('crunchy', 0.9576648473739624),
          ('still', 0.9573863744735718),
          ('edible', 0.9556686878204346)]

In [59]: w2v_model.wv.most_similar('like')

Out[59]: [('taste', 0.8839425444602966),
          ('tastes', 0.8775187134742737),
          ('its', 0.8687935471534729),
          ('smokey', 0.8654364347457886),
          ('sweet', 0.8583647012710571),
          ('smell', 0.8574129939079285),
          ('strong', 0.854999303817749),
          ('doesnt', 0.8540593385696411),
          ('isnt', 0.8467773795127869),
          ('just', 0.8355744481086731)]

In [50]: #count_vect_feat = count_vect.get_feature_names() # list of words in the BoW
         #count_vect_feat.index('like')
         #print(count_vect_feat[64055])

```

## 6 [7.2.7] Avg W2V, TFIDF-W2V

```

In [60]: # 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 list_of_sent: # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        try:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
        except:
            pass
    sent_vec /= cnt_words
    sent_vectors.append(sent_vec)
print(len(sent_vectors))
print(len(sent_vectors[0]))

```

4000  
50

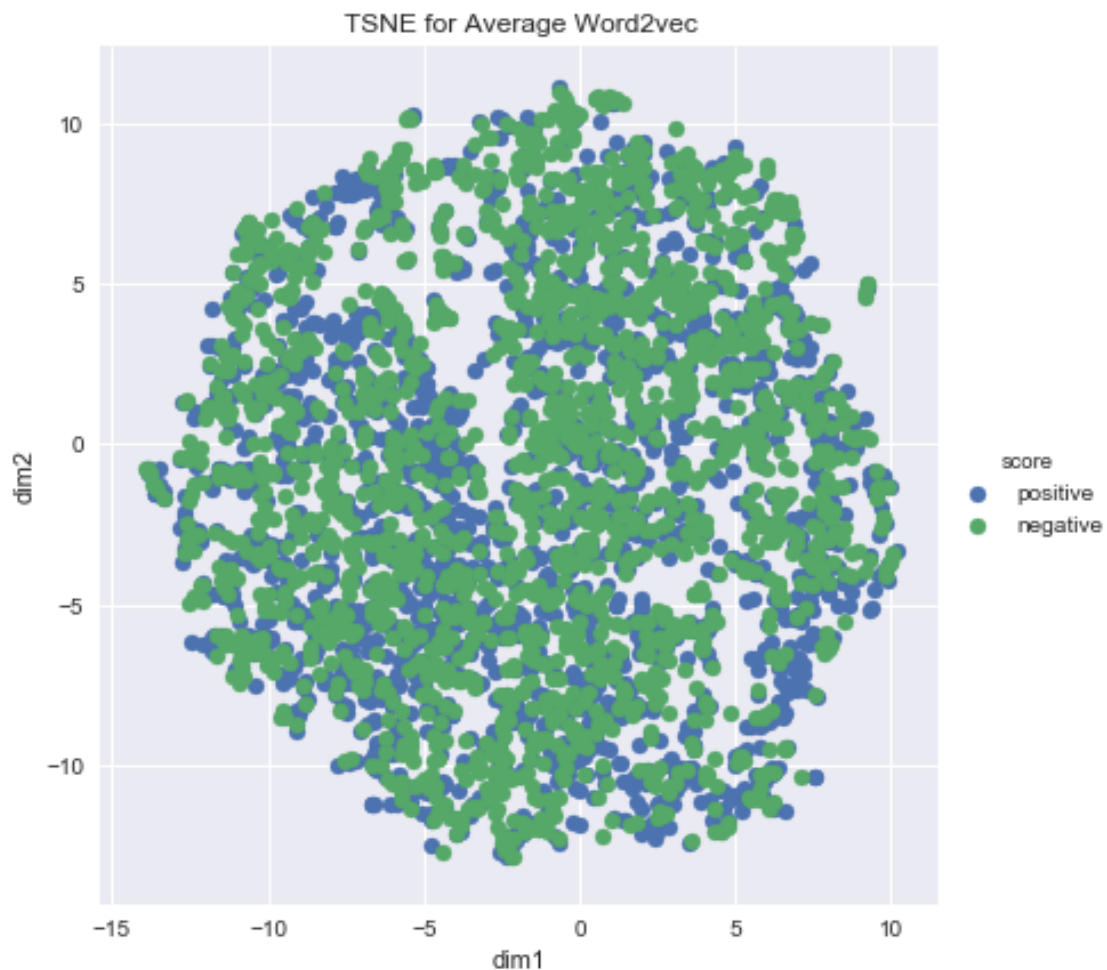
```
In [61]: #sent_vectors = sent_vectors[0:4000]
        #len(sent_vectors)
```

```
In [65]: #tsne
        from sklearn.manifold import TSNE
        model = TSNE(n_components=2, random_state=0, perplexity = 20, n_iter = 5000)

        tsne_data = model.fit_transform(sent_vectors)

        tsne_data = np.vstack((tsne_data.T, score_4000)).T
        tsne_df = pd.DataFrame(data=tsne_data, columns=("dim1", "dim2", "score"))

        # Plotting the result of tsne
        sns.FacetGrid(tsne_df, hue="score", size=6).map(plt.scatter, 'dim1', 'dim2').add_legend()
        plt.title("TSNE for Average Word2vec")
        plt.show()
```



**Observations:-** Here, all +ve and -ve reviews are not well separated this also looks like bow and tfidf vector representations.

```
In [66]: # To avoid warnings
```

```
# http://docs.scipy.org/doc/numpy/reference/generated/numpy.seterr.html
np.seterr(divide='ignore', invalid='ignore')
```

```
Out[66]: {'divide': 'warn', 'invalid': 'warn', 'over': 'warn', 'under': 'ignore'}
```

```
In [67]: # TF-IDF weighted Word2Vec
```

```
tfidf_feat = tf_idf_vect.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 l
row=0;
```

```
for sent in list_of_sent: # 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
```

```
        try:
```

```
            vec = w2v_model.wv[word]
```

```
            # obtain the tf_idfidf of a word in a sentence/review
```

```
            tf_idf = final_tf_idf[row, tfidf_feat.index(word)]
```

```
            sent_vec += (vec * tf_idf)
```

```
            weight_sum += tf_idf
```

```
        except:
```

```
            pass
```

```
    sent_vec /= weight_sum
```

```
    tfidf_sent_vectors.append(sent_vec)
```

```
    row += 1
```

```
In [68]: # To know length of tfidf vector
```

```
len(tfidf_sent_vectors)
```

```
Out[68]: 4000
```

```
In [69]: np.isnan(tfidf_sent_vectors)
```

```
Out[69]: array([[False, False, False, ..., False, False, False],
 [False, False, False, ..., False, False, False],
 [False, False, False, ..., False, False, False],
 ...,
 [False, False, False, ..., False, False, False],
 [False, False, False, ..., False, False, False],
 [False, False, False, ..., False, False, False]], dtype=bool)
```

```
In [71]: # To replace nan with 0 and inf with large finite number
```

```
tfidf_sent_vectors = np.nan_to_num(tfidf_sent_vectors)
```



```

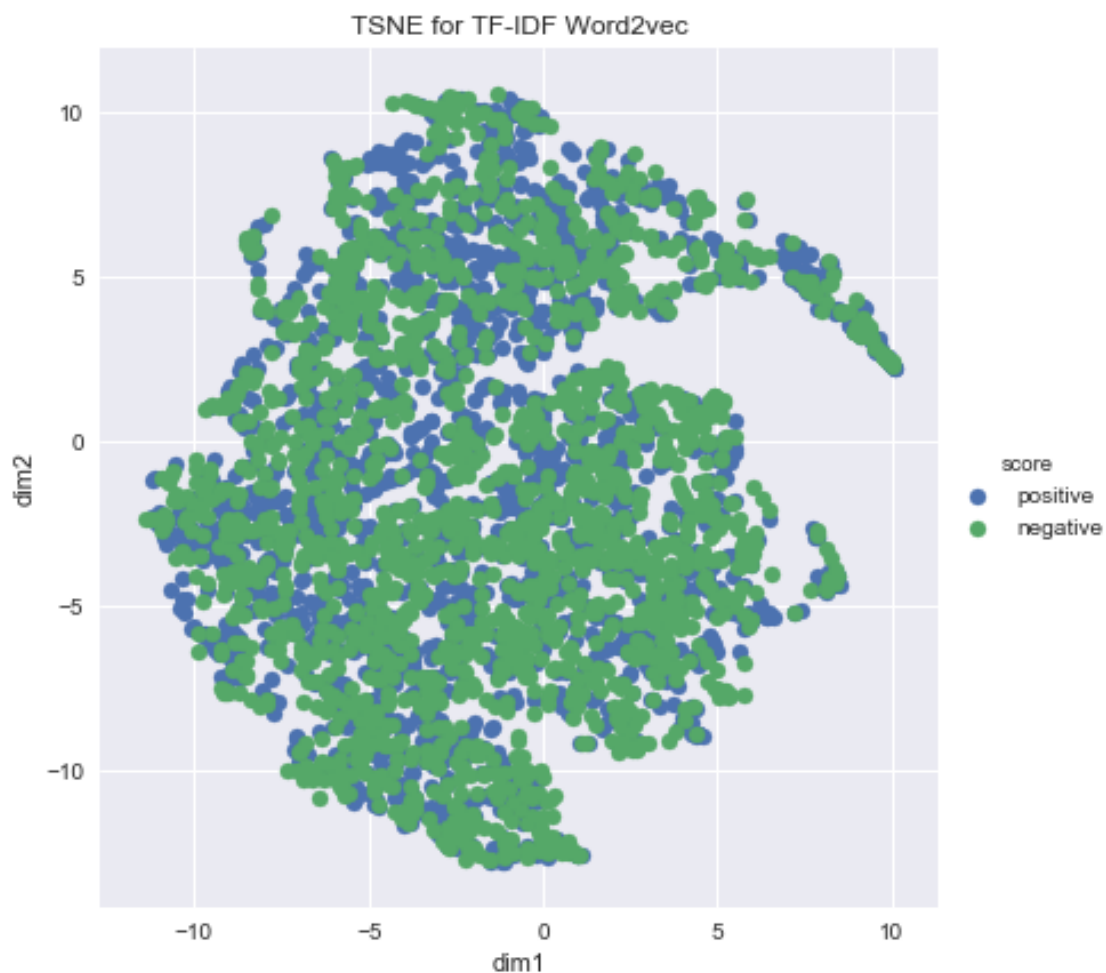
In [72]: #tsne
from sklearn.manifold import TSNE
model = TSNE(n_components=2, random_state=0, perplexity = 50, n_iter = 5000)

tsne_data = model.fit_transform(tfidf_sent_vectors)

tsne_data = np.vstack((tsne_data.T, score_4000)).T
tsne_df = pd.DataFrame(data=tsne_data, columns=("dim1", "dim2", "score"))

# Ploting the result of tsne
sns.FacetGrid(tsne_df, hue="score", size=6).map(plt.scatter, 'dim1', 'dim2').add_legend()
plt.title("TSNE for TF-IDF Word2vec")
plt.show()

```



**Observation-** This plot also looks like the bow, tfidf and avg word2vec. Both +ve and -ve reviews are not well separated they overlapped each other.

**Conclusions:-** 1. As none of TSNE representation gives a well separated both +ve and -ve reviews. 2. We can not simply draw a plane to separate -ve and +ve reviews. Although, By

looking at only visual representation of data we can not take decision whether to draw a plane or not. 3. We will have some alternative method by that we will look at into this problem like how we can separate -ve and +ve reviews.