

Amazon Fine Food Reviews Analysis

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

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

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

Number of reviews: 568,454

Number of users: 256,059

Number of products: 74,258

Timespan: Oct 1999 - Oct 2012

Number of Attributes/Columns in data: 10

Attribute Information:

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

Objective:

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

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

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

[1]. Reading Data

[1.1] Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

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

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

```
In [1]: %matplotlib inline
import warnings
warnings.filterwarnings("ignore")

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
```

```

import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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

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

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

from tqdm import tqdm
import os

```

```

C:\ProgramData\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWarning: detected Windows; aliasing chunkize to chunkize_serial
  warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")

```

```

In [2]: # using 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 50000""", con)

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

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

```

Number of data points in our data (50000, 10)

Out[2]:

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

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

```
In [3]: display = pd.read_sql_query("""
SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
FROM Reviews
GROUP BY UserId
HAVING COUNT(*)>1
""", con)
```

```
In [4]: print(display.shape)
display.head()
```

(80668, 7)

Out[4]:

	UserId	ProductId	ProfileName	Time	Score	Text	COUNT(*)
0	#oc- R115TNMSPFT9I7	B005ZBZLT4	Breyton	1331510400	2	Overall its just OK when considering the price...	2
1	#oc- R11D9D7SHXIJB9	B005HG9ESG	Louis E. Emory "hoppy"	1342396800	5	My wife has recurring extreme muscle spasms, u...	3
2	#oc- R11DNU2NBKQ23Z	B005ZBZLT4	Kim Cieszykowski	1348531200	1	This coffee is horrible and unfortunately not ...	2

	UserId	ProductId	ProfileName	Time	Score	Text	COUNT(*)
3	#oc-R11O5J5ZVQE25C	B005HG9ESG	Penguin Chick	1346889600	5	This will be the bottle that you grab from the...	3
4	#oc-R12KPBODL2B5ZD	B007OSBEV0	Christopher P. Presta	1348617600	1	I didnt like this coffee. Instead of telling y...	2

```
In [5]: display[display['UserId']=='AZY10LLTJ71NX']
```

Out[5]:

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

```
In [6]: display['COUNT(*)'].sum()
```

Out[6]: 393063

[2] Exploratory Data Analysis

[2.1] Data Cleaning: Deduplication

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

```
In [7]: display= pd.read_sql_query("""
```

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

Out[7]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenon
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 it can be seen above that same user has multiple reviews with same values for HelpfulnessNumerator, HelpfulnessDenominator, Score, Time, Summary and Text and on doing analysis it was found that

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

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

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

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

```
In [8]: #Sorting data according to ProductId in ascending order
sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True,
inplace=False, kind='quicksort', na_position='last')
```

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

```
Out[9]: (46072, 10)
```

```
In [10]: #Checking to see how much % of data still remains
(final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

```
Out[10]: 92.144
```

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

too are removed from calculations

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

display.head()
```

```
Out[11]:
```

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenomr
--	----	-----------	--------	-------------	----------------------	-------------------

0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	
1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	

```
In [12]: final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [13]: #Before starting the next phase of preprocessing lets see the number of
entries left
print(final.shape)

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

(46071, 10)
```

```
Out[13]: 1    38479
         0     7592
         Name: Score, dtype: int64
```

[3] Preprocessing

[3.1]. Preprocessing Review Text

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

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

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

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

```
In [14]: # printing some random reviews
sent_0 = final['Text'].values[0]
print(sent_0)
print("="*50)

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

sent_1500 = final['Text'].values[1500]
```

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

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

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

=====

this is yummy, easy and unusual. it makes a quick, delicious pie, crisp or cobbler. home made is better, but a heck of a lot more work. this is great to have on hand for last minute dessert needs where you really want to impress with your creativity in cooking! recommended.

=====

Great flavor, low in calories, high in nutrients, high in protein! Usually protein powders are high priced and high in calories, this one is a great bargain and tastes great, I highly recommend for the lady gym rats, probably not "macho" enough for guys since it is soy based...

=====

For those of you wanting a high-quality, yet affordable green tea, you should definitely give this one a try. Let me first start by saying that everyone is looking for something different for their ideal tea, and I will attempt to briefly highlight what makes this tea attractive to a wide range of tea drinkers (whether you are a beginner or long-time tea enthusiast). I have gone through over 12 boxes of this tea myself, and highly recommend it for the following reasons:

-Quality: First, this tea offers a smooth quality without any harsh or bitter after tones, which often turns people off from many green teas. I've found my ideal brewing time to be between 3-5 minutes, giving you a light but flavorful cup of tea. However, if you get distracted or forget about your tea and leave it brewing for 20+ minutes like I sometimes do, the quality of this tea is such that you still get a smooth but deeper flavor without the bad after taste. The leaves themselves are whole leaves (not powdered stems, branches, etc commonly found in other brands), and the high-quality nylon bags also include chunks of tropical fruit and other discernible ingredients. This isn't your standard cheap paper bag

with a mix of unknown ingredients that have been ground down to a fine powder, leaving you to wonder what it is you are actually drinking.

-Taste: This tea offers notes of real pineapple and other hints of tropical fruits, yet isn't sweet or artificially flavored. You have the foundation of a high-quality young hyson green tea for those true "tea flavor" lovers, yet the subtle hints of fruit make this a truly unique tea that I believe most will enjoy. If you want it sweet, you can add sugar, splenda, etc but this really is not necessary as this tea offers an inherent warmth of flavor through its ingredients.

-Price: This tea offers an excellent product at an exceptional price (especially when purchased at the prices Amazon offers). Compared to other brands which I believe to be of similar quality (Mighty Leaf, Rishi, Two Leaves, etc.), Revolution offers a superior product at an outstanding price. I have been purchasing this through Amazon for less per box than I would be paying at my local grocery store for Lipton, etc.

Overall, this is a wonderful tea that is comparable, and even better than, other teas that are priced much higher. It offers a well-balanced cup of green tea that I believe many will enjoy. In terms of taste, quality, and price, I would argue you won't find a better combination that that offered by Revolution's Tropical Green Tea.

=====

```
In [15]: # remove urls from text python: https://stackoverflow.com/a/40823105/4084039
sent_0 = re.sub(r"http\S+", "", sent_0)
sent_1000 = re.sub(r"http\S+", "", sent_1000)
sent_150 = re.sub(r"http\S+", "", sent_1500)
sent_4900 = re.sub(r"http\S+", "", sent_4900)

print(sent_0)
```

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

```
In [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)

```

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

=====

this is yummy, easy and unusual. it makes a quick, delicious pie, crisp or cobbler. home made is better, but a heck of a lot more work. this is great to have on hand for last minute dessert needs where you really wa nt to impress wih your creativity in cooking! recommended.

=====

Great flavor, low in calories, high in nutrients, high in protein! Usua lly protein powders are high priced and high in calories, this one is a great bargain and tastes great, I highly recommend for the lady gym rat s, probably not "macho" enough for guys since it is soy based...

=====

For those of you wanting a high-quality, yet affordable green tea, you should definitely give this one a try. Let me first start by saying tha t everyone is looking for something different for their ideal tea, and

I will attempt to briefly highlight what makes this tea attractive to a wide range of tea drinkers (whether you are a beginner or long-time tea enthusiast). I have gone through over 12 boxes of this tea myself, and highly recommend it for the following reasons:-Quality: First, this tea offers a smooth quality without any harsh or bitter after tones, which often turns people off from many green teas. I've found my ideal brewing time to be between 3-5 minutes, giving you a light but flavorful cup of tea. However, if you get distracted or forget about your tea and leave it brewing for 20+ minutes like I sometimes do, the quality of this tea is such that you still get a smooth but deeper flavor without the bad after taste. The leaves themselves are whole leaves (not powdered stems, branches, etc commonly found in other brands), and the high-quality nylon bags also include chunks of tropical fruit and other discernible ingredients. This isn't your standard cheap paper bag with a mix of unknown ingredients that have been ground down to a fine powder, leaving you to wonder what it is you are actually drinking.-Taste: This tea offers notes of real pineapple and other hints of tropical fruits, yet isn't sweet or artificially flavored. You have the foundation of a high-quality young hyson green tea for those true "tea flavor" lovers, yet the subtle hints of fruit make this a truly unique tea that I believe most will enjoy. If you want it sweet, you can add sugar, splenda, etc but this really is not necessary as this tea offers an inherent warmth of flavor through its ingredients.-Price: This tea offers an excellent product at an exceptional price (especially when purchased at the prices Amazon offers). Compared to other brands which I believe to be of similar quality (Mighty Leaf, Rishi, Two Leaves, etc.), Revolution offers a superior product at an outstanding price. I have been purchasing this through Amazon for less per box than I would be paying at my local grocery store for Lipton, etc.Overall, this is a wonderful tea that is comparable, and even better than, other teas that are priced much higher. It offers a well-balanced cup of green tea that I believe many will enjoy. In terms of taste, quality, and price, I would argue you won't find a better combination than that offered by Revolution's Tropical Green Tea.

```
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)

```

Great flavor, low in calories, high in nutrients, high in protein! Usually protein powders are high priced and high in calories, this one is a great bargain and tastes great, I highly recommend for the lady gym rats, probably not "macho" enough for guys since it is soy based...

=====

```

In [19]: #remove words with numbers python: https://stackoverflow.com/a/18082370/4084039
sent_0 = re.sub("\S*\d\S*", "", sent_0).strip()
print(sent_0)

```

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

```

In [20]: #remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)

```

```
print(sent_1500)
```

Great flavor low in calories high in nutrients high in protein Usually protein powders are high priced and high in calories this one is a great bargain and tastes great I highly recommend for the lady gym rats probably not macho enough for guys since it is soy based

```
In [21]: # https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
# <br /><br /> ==> after the above steps, we are getting "br br"
# we are including them into stop words list
# instead of <br /> if we have <br> these tags would have been removed in the 1st step

stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \
               "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
               'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their', \
               'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \
               'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
               'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', \
               'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', \
               'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', \
               'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', \
               'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
               's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', \
               've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "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 [22]: # Combining all the above students
from tqdm import tqdm
preprocessed_reviews = []
# tqdm is for printing the status bar
for sentence in tqdm(final['Text'].values):
    sentence = re.sub(r"http\S+", "", sentence)
    sentence = BeautifulSoup(sentence, 'lxml').get_text()
    sentence = decontracted(sentence)
    sentence = re.sub("\S*\d\S*", "", sentence).strip()
    sentence = re.sub('[^A-Za-z]+', ' ', sentence)
    # https://gist.github.com/sebleier/554280
    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower()
() not in stopwords)
    preprocessed_reviews.append(sentence.strip())
```

```
100%|██████████| 46071/46071 [00:35<00:00, 1282.08it/s]
```

```
In [23]: preprocessed_reviews[1500]
```

```
Out[23]: 'great flavor low calories high nutrients high protein usually protein
powders high priced high calories one great bargain tastes great highly
recommend lady gym rats probably not macho enough guys since soy based'
```

```
In [24]: final['CleanedText']=preprocessed_reviews
         final.head(5)
```

Out[24]:

Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
----	-----------	--------	-------------	----------------------	------------------------

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
22620	24750	2734888454	A13ISQV0U9GZIC	Sandikaye	1	
22621	24751	2734888454	A1C298ITT645B6	Hugh G. Pritchard	0	
2546	2774	B00002NCJC	A196AJHU9EASJN	Alex Chaffee	0	
2547	2775	B00002NCJC	A13RRPGE79XFFH	reader48	0	
1145	1244	B00002ZZ754	A3B8RCEI0FXFI6	B G Chase	10	

[3.2] Preprocessing Review Summary

In [25]: `## Similarly you can do preprocessing for review summary also.`

[4] Featurization

[4.1] BAG OF WORDS

```
In [25]: #BoW
count_vect = CountVectorizer() #in scikit-learn
count_vect.fit(preprocessed_reviews)
print("some feature names ", count_vect.get_feature_names()[:10])
print('='*50)

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

some feature names  ['aa', 'aahhs', '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 [26]: #bi-gram, tri-gram and n-gram

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

# you can choose these numebrs min_df=10, max_features=5000, of your choice
count_vect = CountVectorizer(ngram_range=(1,2), min_df=10, max_features=5000)
final_bigram_counts = count_vect.fit_transform(preprocessed_reviews)
```

```
print("the type of count vectorizer ",type(final_bigram_counts))
print("the shape of out text BOW vectorizer ",final_bigram_counts.get_shape())
print("the number of unique words including both unigrams and bigrams "
, final_bigram_counts.get_shape()[1])
```

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

[4.3] TF-IDF

```
In [27]: tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
tf_idf_vect.fit(preprocessed_reviews)
print("some sample features(unique words in the corpus)",tf_idf_vect.get_feature_names()[0:10])
print('='*50)
```

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

```
some sample features(unique words in the corpus) ['ability', 'able', 'able find', 'able get', 'absolute', 'absolutely', 'absolutely delicious', 'absolutely love', 'absolutely no', 'according']
```

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

[4.4] Word2Vec

```
In [28]: # Train your own Word2Vec model using your own text corpus
```

```
i=0
list_of_sentence=[]
for sentence in preprocessed_reviews:
    list_of_sentence.append(sentence.split())
```

In [42]: *# Using Google News Word2Vectors*

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

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

is_your_ram_gt_16g=False
want_to_use_google_w2v = False
want_to_train_w2v = True

if want_to_train_w2v:
    # min_count = 5 considers only words that occurred at least 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', 0.9951335191726685), ('calorie', 0.9946465492248535), ('wonderful', 0.9946032166481018), ('excellent', 0.9944332838058472), ('especially', 0.9941144585609436), ('baked', 0.9940600395202637), ('salted', 0.994047224521637), ('alternative', 0.9937226176261902), ('tasty', 0.9936816692352295), ('healthy', 0.9936649799346924)]
=====
[('varieties', 0.9994194507598877), ('become', 0.9992934465408325), ('popcorn', 0.9992750883102417), ('de', 0.9992610216140747), ('miss', 0.9992451071739197), ('melitta', 0.999218761920929), ('choice', 0.9992102384567261), ('american', 0.9991837739944458), ('beef', 0.9991780519485474), ('finish', 0.9991567134857178)]
```

```
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 Avg W2V, TFIDF-W2V

[4.4.1.1] Avg W2v

```
In [38]: # 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, 1330.47it/s]
```

```
4986
50
```

[4.4.1.2] TFIDF weighted W2v

```
In [39]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer()
tf_idf_matrix = model.fit_transform(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_)))
```

```
100%|███████████| 4986/4986 [00:20<00:00, 245.63it/s]
```

[5] Assignment 11: Truncated SVD

1. **Apply Truncated-SVD on only this feature set:**
- **SET 2:** Review text, preprocessed one converted into vectors using (TFIDF)
 - **Procedure:**

- Take top 2000 or 3000 features from tf-idf vectorizers using `idf_ score`.
- You need to calculate the co-occurrence matrix with the selected features (Note: $X.X^T$ doesn't give the co-occurrence matrix, it returns the covariance matrix, check these blogs [blog-1](#), [blog-2](#) for more information)
- You should choose the `n_components` in truncated svd, with maximum explained variance. Please search on how to choose that and implement them. (hint: plot of cumulative explained variance ratio)
- After you are done with the truncated svd, you can apply K-Means clustering and choose the best number of clusters based on elbow method.
- Print out wordclouds for each cluster, similar to that in previous assignment.
- You need to write a function that takes a word and returns the most similar words using cosine similarity between the vectors(vector: a row in the matrix after truncatedSVD)

Truncated-SVD

[5.1] Taking top features from TFIDF, SET 2

```
In [3]: # Please write all the code with proper documentation
```

```
In [25]: X = final["CleanedText"]
         print("shape of X:", X.shape)
```

```
shape of X: (46071,)
```

```
In [33]: y=X.values
```

```
In [27]: from sklearn.feature_extraction.text import TfidfVectorizer
         vectorizer = TfidfVectorizer(ngram_range=(1,2),min_df=10)
         Xtfidf = vectorizer.fit_transform(X)
         print("the shape of out text TFIDF vectorizer ",Xtfidf.get_shape())
```

the shape of out text TFIDF vectorizer (46071, 27311)

```
In [28]: indice = np.argsort(vectorizer.idf_)[::-1]
feature = vectorizer.get_feature_names()
topn = 2000
topfeatures = [feature[i] for i in indice[:topn]]
```

```
In [30]: topfeat=np.array(topfeatures)
```

```
In [67]: print(topfeat[1:10])
```

```
['left bad' 'course love' 'rooibos tea' 'rose hips' 'leaves bud'
 'rose tea' 'couple cases' 'couple bucks' 'leaves tea']
```

[5.2] Calulation of Co-occurrence matrix

```
In [3]: # Please write all the code with proper documentation
```

```
In [32]: class WordVector:
```

```
    # Initialising the max_features and sample_data to pass in TFIDF vectorizer
```

```
    def __init__(self, max_feat , sample_data):
        self.max_feat = max_feat # No.of top words
        self.sample_data = sample_data # document to vectorize
        # List of all top max_feat words
        self.top_words = []
        self.freq = []
```

```
    # Picking top max_feat words by using TFIDF vextorizer
```

```
    def topWords(self):
        tf_idf_vect = TfidfVectorizer(max_features=self.max_feat)
        tfidf_vec = tf_idf_vect.fit_transform(self.sample_data)
        print("the type of count vectorizer :",type(tfidf_vec))
```

```

print("the shape of out text TFIDF vectorizer : ",tfidf_vec.get_shape())
print("the number of unique words :", tfidf_vec.get_shape()[1])

# Top 'n' words
self.top_words = tf_idf_vect.get_feature_names()
# tfidf frequencies of top 'n' words
self.freq = tf_idf_vect.idf_

return tf_idf_vect.get_feature_names()

# Computing the co-occurrence matrix with value of neighbourhood as neighbour_num
def cooccurrenceMatrix(self, neighbour_num , list_words):

    # Storing all words with their indices in the dictionary
    corpus = dict()
    # List of all words in the corpus
    doc = []
    index = 0
    for sent in self.sample_data:
        for word in sent.split():
            doc.append(word)
            corpus.setdefault(word,[])
            corpus[word].append(index)
            index += 1

    # Co-occurrence matrix
    matrix = []
    # rows in co-occurrence matrix
    for row in list_words:
        # row in co-occurrence matrix
        temp = []
        # column in co-occurrence matrix
        for col in list_words :
            if( col != row):
                # No. of times col word is in neighbourhood of row
                count = 0

```

```

# Value of neighbourhood
num = neighbour_num
# Indices of row word in the corpus
positions = corpus[row]
for i in positions:
    if i < (num-1):
        # Checking for col word in neighbourhood of
row
        if col in doc[i:i+num]:
            count +=1
        elif (i >= (num-1)) and (i <= (len(doc)-num)):
            # Check col word in neighbour of row
            if (col in doc[i-(num-1):i+1]) and (col in
doc[i:i+num]):
                count +=2
            # Check col word in neighbour of row
            elif (col in doc[i-(num-1):i+1]) or (col in
doc[i:i+num]):
                count +=1
        else :
            if (col in doc[i-(num-1):i+1]):
                count +=1

# appending the col count to row of co-occurrence m
atrix
temp.append(count)
else:
    # Append 0 in the column if row and col words are e
qual
    temp.append(0)
# appending the row in co-occurrence matrix
matrix.append(temp)
# Return co-occurrence matrix
return np.array(matrix)

```

In [34]: `wv = WordVector(2000,y)`

```
# Picking top 2K words
words_top = wv.topWords()
```

the type of count vectorizer : <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer : (46071, 2000)
the number of unique words : 2000

```
In [35]: co_occ_matrix = wv.cooccurrenceMatrix(5, words_top)
print("Shape of co-occurrence matrix : ",co_occ_matrix.shape )
print('\n')
```

Shape of co-occurrence matrix : (2000, 2000)

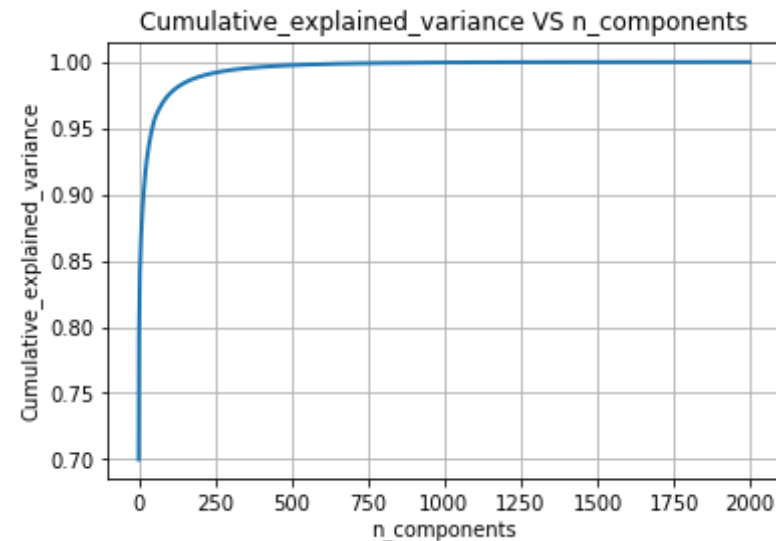
[5.3] Finding optimal value for number of components (n) to be retained.

```
In [3]: # Please write all the code with proper documentation
```

```
In [36]: def plotCV(co_matrix):
    from sklearn.decomposition import TruncatedSVD
    max_features = co_matrix.shape[1]-1
    svd = TruncatedSVD(n_components=max_features)
    svd_data = svd.fit_transform(co_matrix)
    percentage_variance = svd.explained_variance_ / np.sum(svd.explaine
d_variance_)
    cumulative_variance = np.cumsum(percentage_variance)
    # Plot the TruncatedSVD spectrum
    plt.figure(1, figsize=(6, 4))
    plt.clf()
    plt.plot(cumulative_variance, linewidth=2)
    plt.axis('tight')
    plt.grid()
    plt.xlabel('n_components')
    plt.ylabel('Cumulative_explained_variance')
```

```
plt.title("Cumulative_explained_variance VS n_components")
plt.show()
```

```
In [37]: plotCV(co_occ_matrix)
```



Observation: From the plot we can see that after 250 components there is a small increase in percentage so it is better to use only 250 components instead of 2000 components.

[5.4] Applying k-means clustering

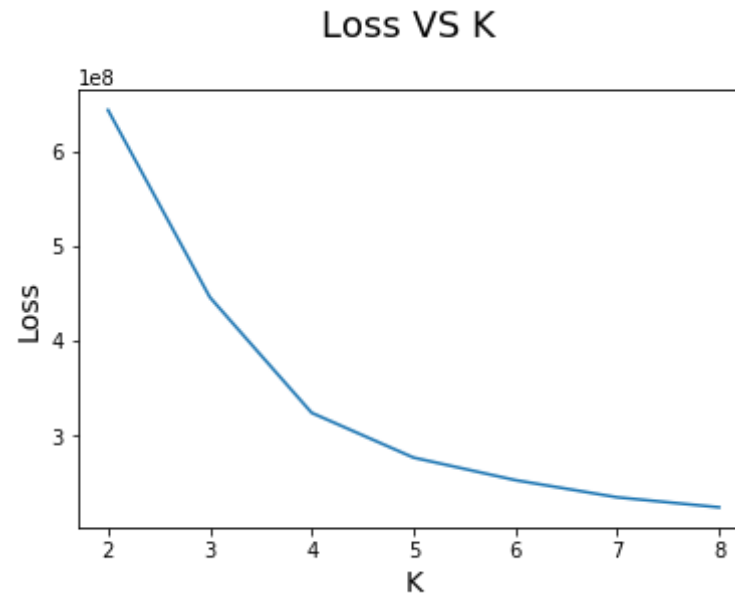
```
In [3]: # Please write all the code with proper documentation
```

```
In [38]: from sklearn.decomposition import TruncatedSVD
svd_trunc = TruncatedSVD(n_components=250)
word_vec = svd_trunc.fit_transform(co_occ_matrix)
print("Shape of word-vector : ", word_vec.shape)
```

```
Shape of word-vector : (2000, 250)
```

```
In [39]: from sklearn.cluster import KMeans
k = [2,3,4,5,6,7,8]
loss = []
for i in k:
    k_means = KMeans(n_clusters=i, n_jobs=-1).fit(word_vec)
    loss.append(k_means.inertia_)
```

```
In [40]: plt.plot(k,loss)
plt.xlabel('K',size=14)
plt.ylabel('Loss',size=14)
plt.title('Loss VS K \n',size=18)
plt.show()
```



```
In [41]: from sklearn.cluster import KMeans
optimalk = 3
k_means = KMeans(n_clusters=optimalk, n_jobs=-1).fit(word_vec)
```

```
In [42]: reviews = final['Text'].values
c1 = []
```

```

c2 = []
c3 = []

for i in range(k_means.labels_.shape[0]):
    if k_means.labels_[i] == 0:
        c1.append(reviews[i])
    elif k_means.labels_[i] == 1:
        c2.append(reviews[i])
    else :
        c3.append(reviews[i])

```

[5.5] Wordclouds of clusters obtained in the above section

In [3]: *# Please write all the code with proper documentation*

```

In [43]: from wordcloud import WordCloud
print("cluster-1")
wordcloud = WordCloud(width = 500, height = 500,
                      background_color = 'white',
                      stopwords = stopwords,
                      min_font_size = 10,max_font_size=40).generate(c1[0])
                      # Display the generated image:
plt.figure(figsize = (8, 8), facecolor = None)
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()

```

cluster-1



```
In [44]: from wordcloud import WordCloud
print("cluster-2")
wordcloud = WordCloud(width = 500, height = 500,
                        background_color = 'white',
                        stopwords = stopwords,
                        min_font_size = 10,max_font_size=40).generate(c2[0])
                        # Display the generated image:
plt.figure(figsize = (8, 8), facecolor = None)
plt.imshow(wordcloud, interpolation='bilinear')
```

```
plt.axis("off")  
plt.show()
```

cluster-2



```
In [45]: from wordcloud import WordCloud  
print("cluster-3")  
wordcloud = WordCloud(width = 500, height = 500,  
                        background_color = 'white',  
                        stopwords = stopwords,  
                        min_font_size = 10,max_font_size=40).generate(c3[0])
```

cluster-3



Create PDF in your applications with the Pdfcrowd [HTML to PDF API](#)

```
In [3]: # Please write all the code with proper documentation
```

```
In [51]: from sklearn.metrics.pairwise import cosine_similarity  
S=cosine_similarity(word_vec[0:1],word_vec)  
S
```

```
Out[51]: array([[1.          , 0.46796461, 0.64881559, ..., 0.61762045, 0.5086350  
2,          0.46519259]])
```

```
In [74]: related_docs_indices = S.argsort()[::-20:-1]  
print(related_docs_indices)  
  
[[ 470 1490  697 ...  983  397    0]]
```

```
In [86]: for i in related_docs_indices:  
         print(topfeat[i])  
  
['two treats' 'cheap chocolate' 'relationship' ... 'stayed fresh'  
'lessons' 'fix without']
```

[6] Conclusions

```
In [1]: # Please write down few lines about what you observed from this assignm  
ent.  
# Also please do mention the optimal values that you obtained for numbe  
r of components & number of clusters.
```

```
In [43]: models = pd.DataFrame({'vectorizer': ["Kmeans with TFIDF"], 'Clusterin  
g' : ["KMeans"], 'Number of cluster': [3]}, columns = ["vectorizer", "Clu  
stering", "Number of cluster"])  
models
```

```
Out[43]:
```

vectorizer	Clustering	Number of cluster
------------	------------	-------------------

	vectorizer	Clustering	Number of cluster
0	Kmeans with TFIDF	KMeans	3

Observation: From the plot we can see that after 250 components there is a small increase in percentage so it is better to use only 250 components instead of 2000 components.

- 1-Text Preprocessing tfidf.
- 2-Finding top 2000 words using TFIDF vectorizer.
- 3-Computing co-occurrence matrix using these top 2000 words.
- 4-Ploted cumulative_explained_variance VS n_components plot to know the correct number of components.
- 5-Applying TruncatedSVD on this co-occurrence matrix with right number of components.
- 6-Applying k-means clustering on this matrix of word-vectors.
- 7-Creating WordClouds for few clusters.
- 8-Cosine similarity of row vector with all other vector.