

Amazon Fine Food Reviews Analysis

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

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

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

Number of reviews: 568,454

Number of users: 256,059

Number of products: 74,258

Timespan: Oct 1999 - Oct 2012

Number of Attributes/Columns in data: 10

Attribute Information:

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

Objective:

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

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

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

[1]. Reading Data

[1.1] Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

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

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

```
In [4]: %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 sklearn.metrics import confusion_matrix, roc_auc_score, roc_curve
from wordcloud import WordCloud, STOPWORDS
from prettytable import PrettyTable
from sklearn.metrics.pairwise import cosine_similarity

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

from tqdm import tqdm
import os

```

```

In [5]: # 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 100000""", 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 (100000, 10)

Out[5]:

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

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

```
In [6]: display = pd.read_sql_query("""
SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
FROM Reviews
GROUP BY UserId
HAVING COUNT(*)>1
""", con)
```

```
In [7]: print(display.shape)
display.head()
```

(80668, 7)

Out[7]:

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

Out[8]:

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

Out[9]: 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 [10]: display= pd.read_sql_query("""
```

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

Out[10]:

	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 [11]: #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 [12]: #Deduplication of entries
final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time",
"Text"}, keep='first', inplace=False)
final.shape
```

```
Out[12]: (87775, 10)
```

```
In [13]: #Checking to see how much % of data still remains
(final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

```
Out[13]: 87.775
```

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

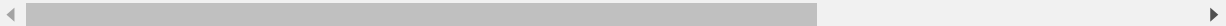
too are removed from calculations

```
In [14]: display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND Id=44737 OR Id=64422
ORDER BY ProductID
""", con)

display.head()
```

```
Out[14]:
```

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



```
In [15]: final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [16]: #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()

(87773, 10)
```

```
Out[16]: 1    73592
```

```
0    14181
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 [17]: from nltk.corpus import stopwords
stop = set(stopwords.words('english')) #set of stopwords
words_to_keep = set(('not'))
stop -= words_to_keep

sno = nltk.stem.SnowballStemmer('english')
def cleanhtml(sentence): #function to clean any HTML Tags
    cleanr = re.compile('<.*?>')
    cleantext = re.sub(cleanr, ' ', sentence)
    return cleantext
```

```
def cleanpunc(sentence): #function to clean any word of punctuation or
special character
    cleaned = re.sub(r'[?!|\\"|\'|#]',r'', sentence)
    cleaned = re.sub(r'[\.,|)|(\|/]',r' ', cleaned)
    return cleaned
```

```
In [18]: #code for implementing step by step check mentioned in preprocessing ph
ase
#runtime will be high due to 500k sentences
i = 0
str1 = ' '
final_string = []
all_positive_words = []
all_negative_words = []
s=' '
for sent in final['Text'].values:
    filtered_sentence=[]
    sent=cleanhtml(sent)
    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)
                    if (final['Score'].values)[i] == 'negative':
                        all_negative_words.append(s)
                else:
                    continue
            else:
                continue
        str1 = b" ".join(filtered_sentence)
        final_string.append(str1)
        i+=1
```

```
In [19]: final['cleanedText']=final_string #Adding a column of Cleanedtext which
displays data after preprocesing.
```

```
final['cleanedText']=final['cleanedText'].str.decode("utf-8")
print('shape of final', final.shape)
final.head()
```

shape of final (87773, 11)

Out[19]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessD
22620	24750	2734888454	A13ISQV0U9GZIC	Sandikaye	1	
22621	24751	2734888454	A1C298ITT645B6	Hugh G. Pritchard	0	
70677	76870	B00002N8SM	A19Q006CSFT011	Arielle	0	
70676	76869	B00002N8SM	A1FYH4S02BW7FN	wonderer	0	
70675	76868	B00002N8SM	AUE8TB5VHS6ZV	eyeofthestorm	0	

In [20]: *# 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 bec ause its a good product but I wont take any chances till they know what is going on with the china imports.

=====

The Candy Blocks were a nice visual for the Lego Birthday party but the candy has little taste to it. Very little of the 2 lbs that I bought w ere eaten and I threw the rest away. I would not buy the candy again.

=====

was way to hot for my blood, took a bite and did a jig lol

=====

My dog LOVES these treats. They tend to have a very strong fish oil sme ll. So if you are afraid of the fishy smell, don't get it. But I think my dog likes it because of the smell. These treats are really small in size. They are great for training. You can give your dog several of the se without worrying about him over eating. Amazon's price was much more reasonable than any other retailer. You can buy a 1 pound bag on Amazon for almost the same price as a 6 ounce bag at other retailers. It's def initely worth it to buy a big bag if your dog eats them a lot.

=====

In [21]: *# 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)
```

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 [22]: *# <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 isn't. It's too bad too because it's a good product but I won't take any chances till they know what is going on with the China imports.

=====

The Candy Blocks were a nice visual for the Lego Birthday party but the candy has little taste to it. Very little of the 2 lbs that I bought were eaten and I threw the rest away. I would not buy the candy again.

=====

was way too hot for my blood, took a bite and did a jig lol

=====

My dog LOVES these treats. They tend to have a very strong fish oil smell. So if you are afraid of the fishy smell, don't get it. But I think my dog likes it because of the smell. These treats are really small in size. They are great for training. You can give your dog several of these without worrying about him over eating. Amazon's price was much more reasonable than any other retailer. You can buy a 1 pound bag on Amazon for almost the same price as a 6 ounce bag at other retailers. It's definitely worth it to buy a big bag if your dog eats them a lot.

```
In [23]: # 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 [24]: sent_1500 = decontracted(sent_1500)
print(sent_1500)
print("="*50)
```

was way to hot for my blood, took a bite and did a jig lol
=====

```
In [25]: #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 bec ause its a good product but I wont take any chances till they know what is going on with the china imports.

```
In [26]: #remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)
print(sent_1500)
```

was way to hot for my blood took a bite and did a jig lol

```
In [27]: # 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', 'yourseve 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 [28]: # 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%|████████████████████████████████████████████████████████████████████████████████| 87773/87773 [01:10<00:00. 1249.76it/s]
```

01/10/2016 10:10:00 AM

'way hot blood took bite jig lol'

[3.2] Preprocessing Review Summary

```
## Similarly you can do preprocessing for review summary also.
# Combining all the above students
from tqdm import tqdm
preprocessed_summaries = []
# 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_summaries.append(sentence.strip())
```

```
| 87773/87773 [01:05<00:00, 1335.96it/s]
```

[4] Featurization

[4.1] BAG OF WORDS

```
#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', 'aaa', 'aaaa', 'aaaaa', 'aaaaaaaaaaaa', 'aaa
aaaaaaaaaaaa', 'aaaaaaahhhhhh', 'aaaaaaaarrrrrrggghhh', 'aaaaaawwwwwwww
w', 'aaaaah']
=====
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer (87773, 54904)
the number of unique words 54904

```

[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 [31]: tf_idf_vect = TfidfVectorizer(ngram_range=(1,1),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) ['aa', 'aafco', 'abac
k', 'abandon', 'abandoned', 'abdominal', 'ability', 'able', 'abroad',
'absence']
```

```
=====
```

```
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer (87773, 11524)
the number of unique words including both unigrams and bigrams 11524
```

```
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.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) ['aa', 'aafco', 'aback', 'abandon', 'abandoned', 'abdominal', 'ability', 'able', 'able add', 'able brew']

=====

the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>

the shape of out text TFIDF vectorizer (87773, 51709)

the number of unique words including both unigrams and bigrams 51709

[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 google'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', '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 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 [27]: # 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 value
dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
```

```
In [41]: # TF-IDF weighted Word2Vec
tfidf_feat = model.get_feature_names() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val = tfidf

tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(list_of_sentence): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum =0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
            # tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole corpus
            # sent.count(word) = tf value of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
```



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

```
In [32]: # Please write all the code with proper documentation
#obtaining features from TfidfVectorizer using idf_score.
idf_score = TfidfVectorizer()
features = tf_idf_vect.get_feature_names()
idf_score = tf_idf_vect.idf_
features = tf_idf_vect.get_feature_names()
idscore_feat = []
for i in range(len(idf_score)):
    idscore_feat.append([idf_score[i], features[i]])
```

```
In [33]: idscore_feat.sort(reverse=True)
idscore_feat=idscore_feat[:3000]
#Top 10 features in idf_feat list
for i in idscore_feat[:10]:
    print(i)
```

```
[9.98462533540777, 'yucca']
[9.98462533540777, 'yogurt']
[9.98462533540777, 'yell']
[9.98462533540777, 'yeasty']
[9.98462533540777, 'yamamotoyama']
[9.98462533540777, 'writeup']
[9.98462533540777, 'wonderful']
[9.98462533540777, 'witch']
[9.98462533540777, 'wiser']
[9.98462533540777, 'winn']
```

[5.2] Calculation of Co-occurrence matrix

```
In [41]: def GetContext(sentence, index):
words = sentence.split(' ')
ret=[]
for word in words:

    if index==0:
        ret.append(words[index+1])
```

```

        ret.append(words[index+2])
    elif index==1:
        ret.append(words[index-1])
        ret.append(words[index+1])
    if len(words)>3:
        ret.append(words[index+2])
    elif index==(len(words)-1):
        ret.append(words[index-2])
        ret.append(words[index-1])
    elif index==(len(words)-2):
        ret.append(words[index-2])
        ret.append(words[index-1])
        ret.append(words[index+1])
    else:
        ret.append(words[index-2])
        ret.append(words[index-1])
        ret.append(words[index+1])
        ret.append(words[index+2])
    return ret

```

```

In [42]: import numpy as np
CORPUS=["abc def ijk pqr", "pqr klm opq", "lmn pqr xyz abc def pqr abc"
]

top2000 = [ "abc","pqr","def"]#list(set((' '.join(ctxs)).split(' ')))
m = np.zeros((3,3))
for sentence in CORPUS:
    for index,word in enumerate(sentence.split(' ')):
        if word in top2000 :
            print(word)
            context=GetContext(sentence,index)
            print(context)
            for word2 in context:
                if word2 in top2000:
                    m[top2000.index(word)][top2000.index(word2)]+=1

print(m)

```

abc

```

'''
['def', 'ijk']
def
['abc', 'ijk']
pqr
['def', 'ijk']
pqr
['klm', 'opq']
pqr
['lmn', 'xyz']
abc
['pqr', 'xyz', 'def', 'pqr']
def
['xyz', 'abc', 'pqr', 'abc']
pqr
['abc', 'def', 'abc']
abc
['def', 'pqr']
[[0. 3. 3.]
 [2. 0. 2.]
 [3. 1. 0.]]
'''

```

```

In [48]: # to obtain the co-occurrence matrix using Top 2000 features of TFIDF v
         ectorizer
         cooccurrenceMatrix = np.zeros((2000,2000)) # co-occurrence matrix
         context_window = 2 # context window for co-occurrence matrix

```

```

In [50]: top_2000_features=[]
         for i in range(2000):
             top_2000_features.append(idfscore_feat[i][1])

```

```

In [51]: len(top_2000_features)

```

```

Out[51]: 2000

```

```

In [52]: # Please write all the code with proper documentation
         for sent in preprocessed_reviews:
             words_sent = sent.split()

```

```

    for index,word in enumerate(words_sent):
        if word in top_2000_features:
            for i in range(max(index - context_window,0),min(index + context_window, len(words_sent) -1) +1):
                if words_sent[i] in top_2000_features:
                    if words_sent[i] != word:
                        cooccurrenceMatrix[top_2000_features.index(words_sent[i]),top_2000_features.index(word)]+=1
                    else:
                        pass
                else:
                    pass
            else:
                pass

```

In [53]: `print(cooccurrenceMatrix)`

```

[[0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]
 ...
 [0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]]

```

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

In [56]:

```

# Please write all the code with proper documentation
# Program to find the optimal number of components for Truncated SVD
from sklearn.decomposition import TruncatedSVD
n_comp = [1,4,10,15,20,50,100,150,200,500,700,800,900,1000,1500, 2000]
# list containing different values of components
explained = [] # explained variance ratio for each component of Truncated SVD
for x in n_comp:
    svd = TruncatedSVD(n_components=x)

```

```

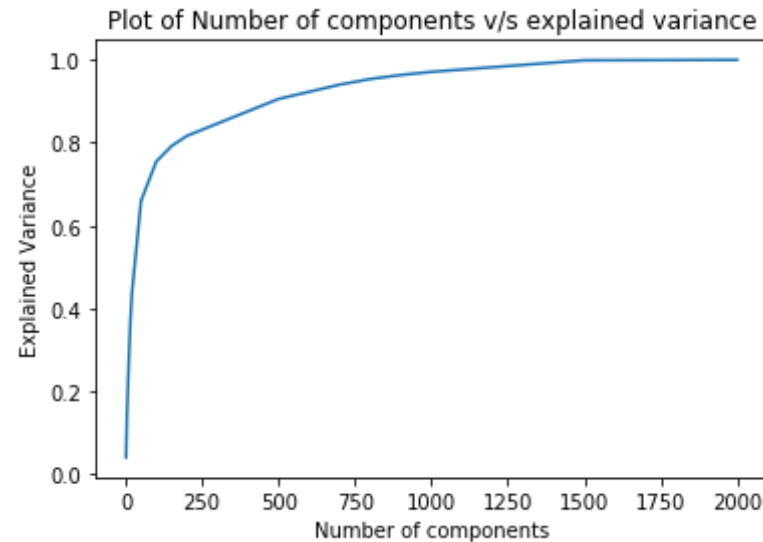
svd.fit(cooccurrenceMatrix)
explained.append(svd.explained_variance_ratio_.sum())
print("Number of components = %r and explained variance = %r"%(x,sv
d.explained_variance_ratio_.sum()))
plt.plot(n_comp, explained)
plt.xlabel('Number of components')
plt.ylabel("Explained Variance")
plt.title("Plot of Number of components v/s explained variance")
plt.show()

```

```

Number of components = 1 and explained variance = 0.04119403781718986
Number of components = 4 and explained variance = 0.14255355508662088
Number of components = 10 and explained variance = 0.27319150611505766
Number of components = 15 and explained variance = 0.36728669832258742
Number of components = 20 and explained variance = 0.43472370967170298
Number of components = 50 and explained variance = 0.65972393922387529
Number of components = 100 and explained variance = 0.75487405821182119
Number of components = 150 and explained variance = 0.79215187201458759
Number of components = 200 and explained variance = 0.81655120104402323
Number of components = 500 and explained variance = 0.90569068406677888
Number of components = 700 and explained variance = 0.94018651466613745
Number of components = 800 and explained variance = 0.95372829680184701
Number of components = 900 and explained variance = 0.96344399041018569
Number of components = 1000 and explained variance = 0.9710124290257078
3
Number of components = 1500 and explained variance = 0.9989071223026705
4
Number of components = 2000 and explained variance = 1.00000000000000078

```



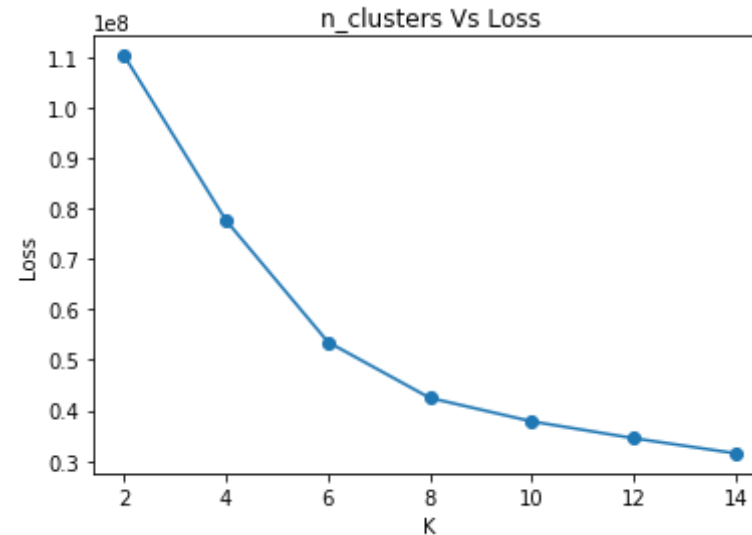
Observation

- From the above plot we can infer that at 500 components having almost 99% variance explained.
- Taking `n_components` a 500 instead of 2000 components.
- Again fit the model with train data.

[5.4] Applying k-means clustering

```
In [52]: # Please write all the code with proper documentation
from sklearn.cluster import KMeans
#Applying optimal n_components to TruncatedSVD to find our train data
tsvd = TruncatedSVD(n_components=500)
X_train = tsvd.fit_transform(cooccurrenceMatrix)
```

```
#Elbow method to find optimal K
def find_optimal_k(data):
    loss = []
    k = list(range(2, 15, 2))
    for noc in k:
        model = KMeans(n_clusters = noc)
        model.fit(data)
        loss.append(model.inertia_)
    plt.plot(k , loss, "-o")
    plt.title('n_clusters Vs Loss')
    plt.xlabel('K')
    plt.ylabel('Loss')
    plt.show()
# Find best K using elbow method
find_optimal_k(X_train)
```



```
In [54]: #training model with optimal k=6
model = KMeans(n_clusters=6).fit(X_train)
```

Observation

- Using KMeans found optimal K=6 by implementing elbow method
- where the point of inflection can be taken as best measure.
- plotted word cloud for each of the clusters.
- Using TSVD train the model with n_comp we got.
- fit the model accordingly.

[5.5] Wordclouds of clusters obtained in the above section

```
In [55]: # Please write all the code with proper documentation
cluster1, cluster2, cluster3, cluster4, cluster5 = [], [], [], [], []
for i in range(model.labels_.shape[0]):
    if model.labels_[i] == 0:
        cluster1.append(top_3000_features[i])
    elif model.labels_[i] == 1:
        cluster2.append(top_3000_features[i])
    elif model.labels_[i] == 2:
        cluster2.append(top_3000_features[i])
    elif model.labels_[i] == 3:
        cluster3.append(top_3000_features[i])
    elif model.labels_[i] == 4:
        cluster4.append(top_3000_features[i])
    else:
        cluster5.append(top_3000_features[i])
```

```
In [59]: # Please write all the code with proper documentation
#for cluster 1
data=''
for i in cluster1:
    data=data+i+' '
from wordcloud import WordCloud
wordcloud = WordCloud(background_color="black").generate(data)

# Display the generated image:
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```



```
In [61]: # Please write all the code with proper documentation
#cluster3
data=''
for i in cluster3:
    data=data+i+' '
from wordcloud import WordCloud
wordcloud = WordCloud(background_color="black").generate(data)

# Display the generated image:
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```



```
In [62]: # Please write all the code with proper documentation
#cluster4
data=''
for i in cluster4:
    data=data+i+' '
from wordcloud import WordCloud
wordcloud = WordCloud(background_color="black").generate(data)

# Display the generated image:
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```

dark

```
In [63]: # Please write all the code with proper documentation
#cluster5
data=''
for i in cluster5:
    data=data+i+' '
from wordcloud import WordCloud
wordcloud = WordCloud(background_color="black").generate(data)

# Display the generated image:
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```



[5.6] Function that returns most similar words for a given word.

```
In [59]: # Please write all the code with proper documentation
from sklearn.metrics.pairwise import cosine_similarity
def similar_word(word, n):
    top_words = []
    similarity = cosine_similarity(cooccurrenceMatrix)
    word_vect = similarity[top_2000_features.index(word)]
    index = word_vect.argsort()[::-1][1:11]
    for i in range(len(index)):
        print((i+1), "Word", top_2000_features[index[i]], "is similar to", word, "\n")
```

```
In [62]: print('Top 10 words similar to witch: ')
print(similar_word('witch', 10))
```

```
Top 10 words similar to witch:
1 Word remedied is similar to witch

2 Word ruby is similar to witch

3 Word rubbish is similar to witch

4 Word romano is similar to witch

5 Word romaine is similar to witch

6 Word rodents is similar to witch

7 Word risky is similar to witch

8 Word reveal is similar to witch

9 Word responding is similar to witch

10 Word resident is similar to witch

None
```

```
In [67]: print('Top 10 words similar to wonderful: ')
print(similar_word('wonderful', 10))
```

```
Top 10 words similar to wonderful:
1 Word wonderful is similar to wonderful

2 Word hollywood is similar to wonderful

3 Word brooklyn is similar to wonderful

4 Word ruth is similar to wonderful

5 Word rodents is similar to wonderful

6 Word risky is similar to wonderful

7 Word romaine is similar to wonderful

8 Word reveal is similar to wonderful

9 Word renew is similar to wonderful

10 Word romano is similar to wonderful

None
```

[6] Conclusions

Please write down few lines about what you observed from this assignment.

Also please do mention the optimal values that you obtained for number of components & number of clusters.

- Considered 100k data points from the whole data set.
- Using idf_score as parameter for TFIDF Vectorizer, find n_components.
- Here we got n_components = 500, so using these components find optimal KMeans.

- find best K using KMeans, optimal K=6 fit the model.
- plotted WordCloud for each of the clusters we got.
- Observed Different Word Cluster for each of the clusters implemented.
- Defined Cosine Similarity Function to get the most similar words for a given word.
- printed top 10 words similar to word given.