

Amazon Fine Food Reviews Analysis

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

EDA: <https://nycdatascience.com/blog/student-works/amazon-fine-foods-visualization/>
(<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 wil be set to "positive". Otherwise, it will be set to "negative".

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

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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
```

```
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 500000 data points  
# you can change the number to any other number based on your computing power  
  
# filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000""", con)  
# for tsne assignment you can take 5k data points  
  
filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3 LIMIT 10000""", 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[2]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0
2	3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1

In [3]:

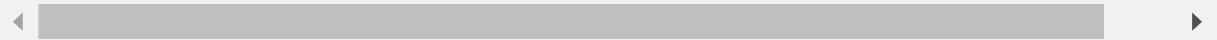
```
display = pd.read_sql_query("""
SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
FROM Reviews
GROUP BY UserId
HAVING COUNT(*)>1
""", con)
```

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

(80668, 7)

Out[4]:

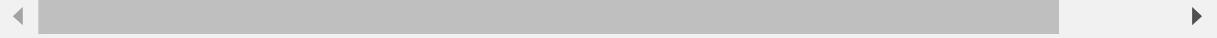
	UserId	ProductId	ProfileName	Time	Score	Text	COU
0	#oc-R115TNMSPFT9I7	B007Y59HVM	Breyton	1331510400	2	Overall its just OK when considering the price...	2
1	#oc-R11D9D7SHXIJB9	B005HG9ET0	Louis E. Emory "hoppy"	1342396800	5	My wife has recurring extreme muscle spasms, u...	3
2	#oc-R11DNU2NBKQ23Z	B007Y59HVM	Kim Cieszykowski	1348531200	1	This coffee is horrible and unfortunately not ...	2
3	#oc-R11O5J5ZVQE25C	B005HG9ET0	Penguin Chick	1346889600	5	This will be the bottle that you grab from the...	3
4	#oc-R12KPBODL2B5ZD	B007OSBE1U	Christopher P. Presta	1348617600	1	I didnt like this coffee. Instead of telling y...	2



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

Out[5]:

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



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

```
Out[6]: 393063
```

[2] Exploratory Data Analysis

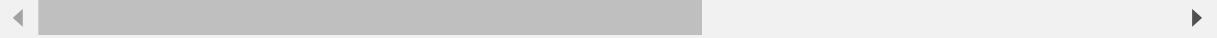
[2.1] Data Cleaning: Deduplication

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	HelpfulnessDenominator
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	2
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	2
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	2
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	2
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan	2	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]: (87775, 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]: 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 [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	Helpful
0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	1
1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	2

◀ ▶

```
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()
```

```
(87773, 10)
```

Out[13]:

```
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 [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 isn't. It's 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.

=====

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 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 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 [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 isn't. It's 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 isn't. It's 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.

=====

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 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 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 [17]: [# https://stackoverflow.com/a/47091490/4084039](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)
```

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

=====

In [19]: [#remove words with numbers python: https://stackoverflow.com/a/18082370/4084039](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 isn't. It's 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](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 [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 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', 'mo
re', \
    'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
    's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'how', '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', "isn't", 'ma', 'mightn', "mightn't", 'mustn', \
    "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'were
n', "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% [██████████] 87773/87773 [00:41<00:00, 2124.28it/s]

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

Out[23]: 'way hot blood took bite jig lol'

[4] Featurization

[4.1] BAG OF WORDS

In [24]: #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', 'aaaaaaaaaaaaaaaaaa', 'aaaaaaaaaaaaaa
hhhhhh', 'aaaaaaaaaaaaaaaaaaaaa', 'aaaaaaaaaaaaaaaaaaaaaa', 'aaaaaaaaaaaaaaaaaaaaaa
=====
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 [25]: #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 numbers 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 (87773, 5000)
the number of unique words including both unigrams and bigrams 5000
```

[4.3] TF-IDF

```
In [26]: tf_idf_vect = TfidfVectorizer(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) ['aa', 'aafco', 'aback', '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

[4.4] Word2Vec

```
In [27]: # Train your own Word2Vec model using your own text corpus
i=0
list_of_sentance=[]
for sentance in preprocessed_reviews:
    list_of_sentance.append(sentance.split())
```

In [28]: # 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 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/0B7XkCwpI5KDYNlNUTbSISS21pQmM/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 atleast 5 times
    w2v_model=Word2Vec(list_of_senteance,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 ")
```

```
[('fantastic', 0.8474854230880737), ('terrific', 0.8286346793174744), ('good', 0.819901704788208), ('awesome', 0.8162736296653748), ('excellent', 0.8037559390068054), ('wonderful', 0.741725504398346), ('perfect', 0.7410205602645874), ('amazing', 0.7299259305000305), ('decent', 0.6918753385543823), ('nice', 0.6916672587394714)]
=====
[('greatest', 0.7970951199531555), ('best', 0.7072246074676514), ('tastiest', 0.7056244611740112), ('horrible', 0.6568552851676941), ('cry', 0.6518000960350037), ('nastiest', 0.6455205678939819), ('disgusting', 0.6337440013885498), ('awful', 0.6105102896690369), ('smoothest', 0.6097800731658936), ('nicest', 0.6007869243621826)]
```

```
In [29]: 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 17386
sample words ['thatthe', 'whisk', 'viva', 'sharper', 'sweet', 'devout', 'cos', 'bash', 'granuals', 'yukky', 'theme', 'bay', 'quibble', 'sanitary', 'days', 'downing', 'mandatory', 'previous', 'diners', 'bride', 'chase s', 'jamaican', 'caffe', 'veterinarians', 'vegetative', 'preoccupied', 'dunked', 'dismissed', 'lawry', 'ruini ng', 'solidified', 'satifying', 'dd', 'chomp', 'different', 'tractor', 'lucaffe', 'pigskin', 'saut', 'callebaut', 'nu tritionists', 'bosch', 'ooh', 'draws', 'economy', 'unremarkable', 'ambient', 'melted', 'comes', 'romanof f']

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

[4.4.1.1] Avg W2v

```
In [30]: # average Word2Vec
# compute average word2vec for each review.
sent_vectors = [] # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(list_of_sentance): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to change th
is to 300 if you use google's w2v
    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
        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% [██████████] 87773/87773 [18:34<00:00, 78.73it/s]

87773

50

[4.4.1.2] TFIDF weighted W2v

```
In [31]: # 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 [32]: # 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_senteance): # 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 courpus
            # sent.count(word) = tf valeus of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    tfidf_sent_vectors.append(sent_vec)
    row += 1
```

100% [██████████] 87773/87773 [1:31:07<00:00, 16.05it/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: $X \cdot X^T$ doesn't give the co-occurrence matrix, it returns the covariance matrix, check these blogs [blog-1](https://medium.com/data-science-group-iitr/word-embedding-2d05d270b285), (<https://medium.com/data-science-group-iitr/word-embedding-2d05d270b285>) [blog-2](https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/) (<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/>) 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 [33]: #idf score
idf_score = tf_idf_vect.idf_
#feature names
ft = tf_idf_vect.get_feature_names()
print(idf_score.shape)
print(len(ft))
```

```
(11524,)
11524
```

```
In [125]: df = pd.DataFrame(list(zip(idf_score,ft)), columns=['idf_score', 'features'])
df = df.sort_values(by=['idf_score'], ascending=False)
top3k_ft = list(df.features[:3000])
len(top3k_ft)
```

```
Out[125]: 3000
```

[5.2] Calculation of Co-occurrence matrix

```
In [126]: #initializing cooccurrence matrix to 0
co_occur=np.zeros((3000,3000))
#taking context_window=2
context_window=2
```

```
In [127]: %%time
for sentance in preprocessed_reviews:
    words=sentance.split()
    for idx,word in enumerate(words):
        if word in top3k_ft:
            for i in range(max(0,idx-context_window),min(idx+context_window,len(words))):

                if words[i] in top3k_ft and words[i]!=word:
                    co_occur[top3k_ft.index(words[i]),top3k_ft.index(word)]+=1
```

Wall time: 3min 12s

```
In [128]: #cooccurrence matrix
print(co_occur)
```

```
[[ 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 [129]: **%%time**

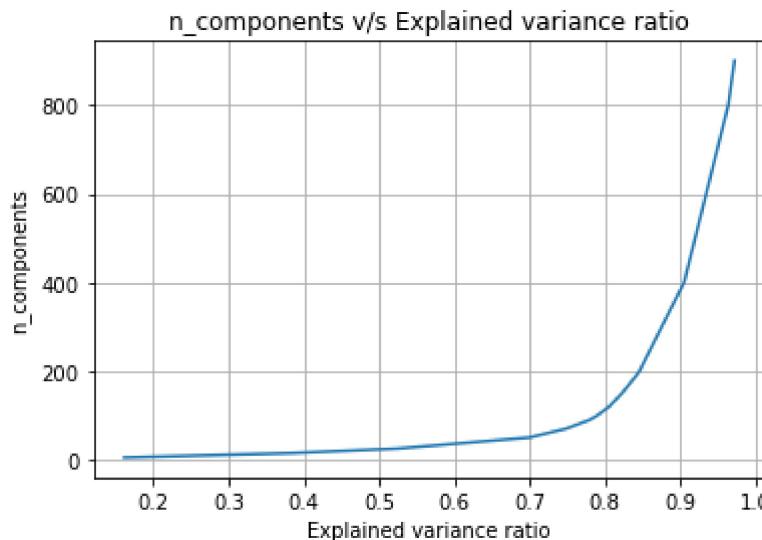
```
from sklearn.decomposition import TruncatedSVD
K = [5,15,25,50,70,90,100,120,150,190,200,400,800,900]
exp_var=[]

for k in K:
    trunc = TruncatedSVD(n_components = k)
    trunc.fit(co_occur)
    exp = trunc.explained_variance_ratio_.sum()
    exp_var.append(exp)
    print('n_components=',k,'Explained variance ratio=',exp)
```

```
n_components= 5 Explained variance ratio= 0.161310544382
n_components= 15 Explained variance ratio= 0.38588941089
n_components= 25 Explained variance ratio= 0.523758079745
n_components= 50 Explained variance ratio= 0.699630053983
n_components= 70 Explained variance ratio= 0.748557016072
n_components= 90 Explained variance ratio= 0.780050358093
n_components= 100 Explained variance ratio= 0.79008544346
n_components= 120 Explained variance ratio= 0.805712400337
n_components= 150 Explained variance ratio= 0.822777397379
n_components= 190 Explained variance ratio= 0.841801570101
n_components= 200 Explained variance ratio= 0.846001714911
n_components= 400 Explained variance ratio= 0.905240365164
n_components= 800 Explained variance ratio= 0.963949592979
n_components= 900 Explained variance ratio= 0.972180123924
Wall time: 17.4 s
```

In [130]: *#plotting curve between n_components and explained variance*

```
plt.plot(exp_var,K)
plt.ylabel('n_components')
plt.xlabel("Explained variance ratio")
plt.title("n_components v/s Explained variance ratio")
plt.grid()
plt.show()
```



[5.4] Applying k-means clustering

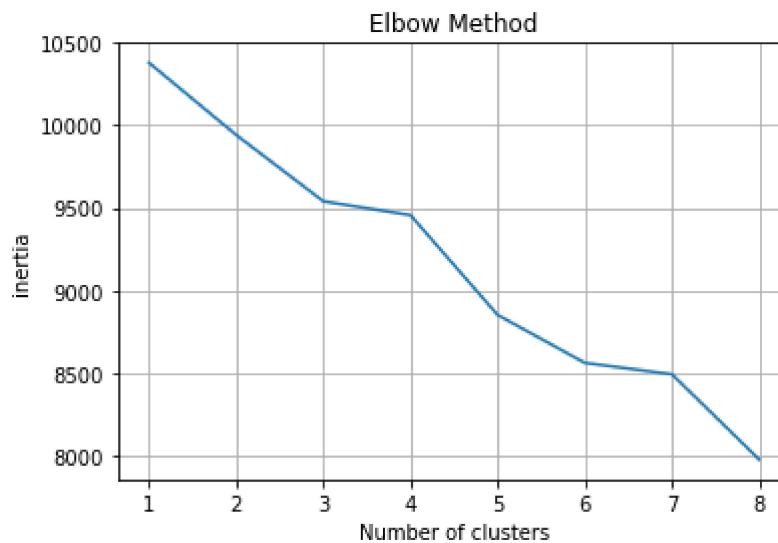
applying truncated svd with best n_components

```
In [131]: from sklearn.cluster import KMeans
trunc2 = TruncatedSVD(n_components=200)
final_data = trunc2.fit_transform(co_occur)
```

finding the best k cluster for model

```
In [132]: %%time
from sklearn.cluster import KMeans
dist = []
# Iterate from 1-9
for i in range(1,9):
    # Initialize KMeans algorithm
    km=KMeans(n_clusters=i,init='k-means++', random_state=0)
    # Fit on data
    km.fit(final_data)
    dist.append(km.inertia_)
# plt.figure(figsize=(10,10))

plt.plot(range(1,9),dist)
plt.title('Elbow Method')
plt.xlabel('Number of clusters')
plt.ylabel('inertia')
plt.grid()
plt.show()
```



Wall time: 2 s

modeling with best number of k cluster

```
In [133]: %%time
#best k= 3
final_model=KMeans(n_clusters=3,precompute_distances=True)
final_model.fit(final_data)
```

Wall time: 209 ms

Visualising the cluster

```
In [134]: labels = list(final_model.labels_)
tt = top3k_ft
new_df = pd.DataFrame(list(zip(tt, labels)),columns =['review', 'cluster'])
cluster_1 = new_df[new_df['cluster']==0]
cluster_2 = new_df[new_df['cluster']==1]
cluster_3 = new_df[new_df['cluster']==2]
```

[5.5] Wordclouds of clusters obtained in the above section

function for wordcloud

```
In [135]: def wordcloud(review):
    string=""
    for i in review:
        string=string+i+' '
    from wordcloud import WordCloud
    wordcloud = WordCloud(background_color="red").generate(string)
    plt.figure(figsize=(10, 8))
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis("off")
    plt.show()
```

In [136]: **%%time**

```
print('WordCloud of cluster 1 ')
wordcloud(cluster_1.review)
```

WordCloud of cluster 1



Wall time: 572 ms

In [137]: **print('WordCloud of cluster 2 ')**
wordcloud(cluster_2.review)

WordCloud of cluster 2



```
In [138]: print('WordCloud of cluster 3 ')
wordcloud(cluster_3.review)
```

WordCloud of cluster 3



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

```
In [149]: #https://stackoverflow.com/questions/29484529/cosine-similarity-between-two-words-in-a-list
#https://www.machinelearningplus.com/nlp/cosine-similarity/
from sklearn.metrics.pairwise import cosine_similarity
def sim_words(word):
    top=[]
    sim=cosine_similarity(co_occur)
    val=sim[top3k_ft.index(word)]
    index=np.argsort(val)
    for i in range(1,11):
        top.append(top3k_ft[index[i]])
    print(top)
```

```
In [150]: #printing top 10 similar words
sim_words('sciences')
```

```
['folk', 'taiwan', 'poly', 'slab', 'dreaming', 'dreamed', 'blanch', 'hybrid', 'twining', 'reliably']
```

[6] Conclusions

1. We have taken top 3000 features of tfidf vectors using tfidf score
2. Then we found co co urence matrix of those feature with context window of 2
3. With the help of explained variance ratio plot, we have found the optimal K for truncated SVD
4. We have taken n_components(k) = 200 as it explained 84.6 % of variance
5. Then applying KMeans cluster on top of it by using elbow method , finding the optimal K cluster = 3
6. We have seen some clusters are dense some have only few words then visualising with word clouds
7. In the end created a function which returns similar words from our training dataset using cosine similarity

In []: