

# 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



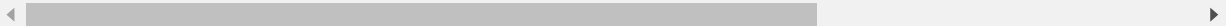
too are removed from calculations

```
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	HelpfulnessDenominator
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:<br /><br />-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.<br /><br />-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.<br /><br />-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.<br /><br />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(r"\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(r'[\^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', "isn'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 [00:33<00:00, 1382.83it/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 [6]: `## 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]
```



- **SET 3:** Review text, preprocessed one converted into vectors using (AVG W2v)
- **SET 4:** Review text, preprocessed one converted into vectors using (TFIDF W2v)
- Find the best 'k' using the elbow-knee method (plot k vs inertia\_)
- Once after you find the k clusters, plot the word cloud per each cluster so that at a single go we can analyze the words in a cluster.

## 2. Apply Agglomerative Clustering on these feature sets:

- **SET 3:** Review text, preprocessed one converted into vectors using (AVG W2v)
- **SET 4:** Review text, preprocessed one converted into vectors using (TFIDF W2v)
- Apply agglomerative algorithm and try a different number of clusters like 2,5 etc.
- Same as that of K-means, plot word clouds for each cluster and summarize in your own words what that cluster is representing.
- You can take around 5000 reviews or so (as this is very computationally expensive one)

## 3. Apply DBSCAN Clustering on these feature sets:

- **SET 3:** Review text, preprocessed one converted into vectors using (AVG W2v)
- **SET 4:** Review text, preprocessed one converted into vectors using (TFIDF W2v)
- Find the best 'Eps' using the [elbow-knee method](#).
- Same as before, plot word clouds for each cluster and summarize in your own words what that cluster is representing.
- You can take around 5000 reviews for this as well.

# [5.1] K-Means Clustering

## [5.1.1] Applying K-Means Clustering on BOW, **SET 1**

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

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

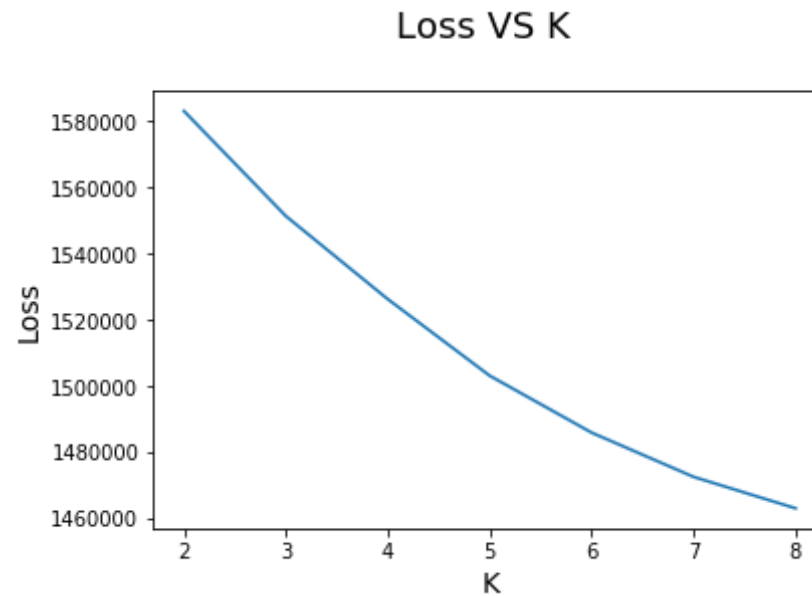
shape of X: (46071,)
```

```
In [26]: #Bow
count_vect = CountVectorizer(min_df = 10,max_features=500) #in scikit-learn
final_counts = count_vect.fit_transform(X)
print("the type of count vectorizer ",type(final_counts))
print("the shape of out text BOW vectorizer ",final_counts.get_shape())

the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer (46071, 500)
```

```
In [27]: 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(final_counts)
    loss.append(k_means.inertia_)
```

```
In [28]: 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 [28]: from sklearn.cluster import KMeans
bow_optimalk = 4
k_means = KMeans(n_clusters=bow_optimalk, n_jobs=-1).fit(final_counts)
```

```
In [30]: reviews = final['Text'].values
c1 = []
c2 = []
c3 = []
c4 = []

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])
    elif k_means.labels_[i] == 2:
        c3.append(reviews[i])
    else :
```

```
c4.append(reviews[i])
```

### [5.1.2] Wordclouds of clusters obtained after applying k-means on BOW **SET 1**

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

```
In [31]: from wordcloud import WordCloud
print("cluster-1")
for i in range(2):
    print("review",i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10,max_font_size=40).generate(c1[i])
    # Display the generated image:
    plt.figure(figsize = (8, 8), facecolor = None)
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis("off")
    plt.show()
```

```
cluster-1
review 0
```

A word cloud visualization of text data. The words are arranged in a circular pattern, with their size and color indicating their frequency or importance. The words include: "isnt", "know", "imports", "good", "till", "find", "wont", "anymore", "made", "going", "product", "USA", "loves", "chicken", "one", "buying", "China", "chances", "hard", "take", "bad", and "dogs". The words are colored in shades of green, blue, and yellow.

review 1



Observation: Cluster-1 is about pet food quality in different countries.

```
In [32]: from wordcloud import WordCloud
print("cluster-2")
for i in range(2):
    print("review",i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10,max_font_size=40).generate(c2[i])
```

```
# Display the generated image:  
plt.figure(figsize = (8, 8), facecolor = None)  
plt.imshow(wordcloud, interpolation='bilinear')  
plt.axis("off")  
plt.show()
```

cluster-2  
review 0



review 1

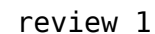


Observation:Cluster-2 is about drinks like tea,freash juice etc

```
In [33]: from wordcloud import WordCloud
print("cluster-3")
for i in range(2):
    print("review",1)
    wordcloud = WordCloud(width = 500, height = 500,
        background_color = 'white',
        stopwords = stopwords,
        min_font_size = 10,max_font_size=40).generate(c3[i])
```

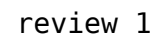


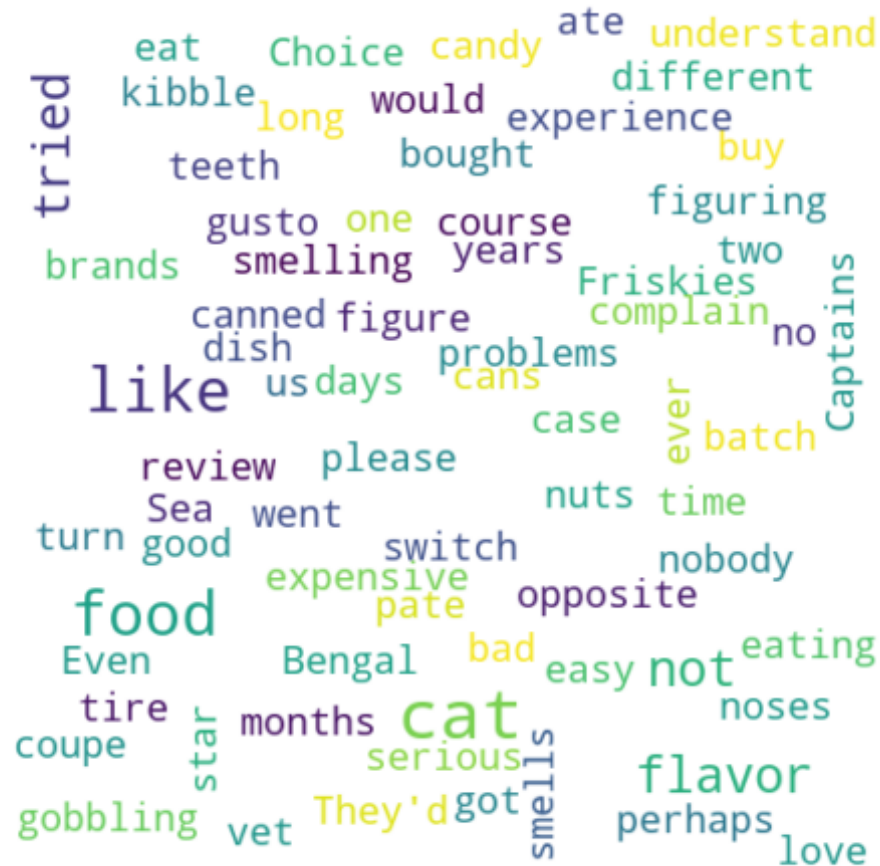
cluster-3  
review 1





cluster-4  
review 0





Observation: Cluster-4 is about pet foods for dogs, cats, fish etc

### [5.1.3] Applying K-Means Clustering on TFIDF, SET 2

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

```
In [35]: X = final["CleanedText"]  
print("shape of X:", X.shape)
```

shape of X: (46071,)

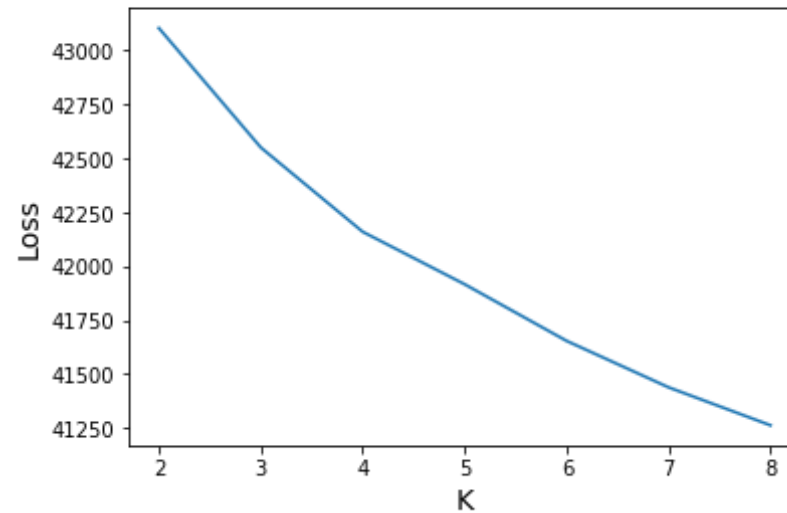
```
In [36]: tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10,max_features=500)
final_tf_idf = tf_idf_vect.fit_transform(X)
print("the type of count vectorizer ",type(final_tf_idf))
print("the shape of out text TFIDF vectorizer ",final_tf_idf.get_shape())
```

the type of count vectorizer <class 'scipy.sparse.csr.csr\_matrix'>  
the shape of out text TFIDF vectorizer (46071, 500)

```
In [28]: 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(final_tf_idf)
    loss.append(k_means.inertia_)
```

```
In [29]: plt.plot(k,loss)
plt.xlabel('K',size=14)
plt.ylabel('Loss',size=14)
plt.title('Loss VS K \n',size=18)
plt.show()
```

Loss VS K



```
In [37]: tfidf_optimalk = 4
k_means = KMeans(n_clusters=tfidf_optimalk, n_jobs=-1).fit(final_tf_idf
)
```

```
In [38]: reviews = final['Text'].values
c1 = []
c2 = []
c3 = []
c4 = []

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])
    elif k_means.labels_[i] == 2:
        c3.append(reviews[i])
    else :
```

```
c4.append(reviews[i])
```

#### [5.1.4] Wordclouds of clusters obtained after applying k-means on TFIDF SET 2

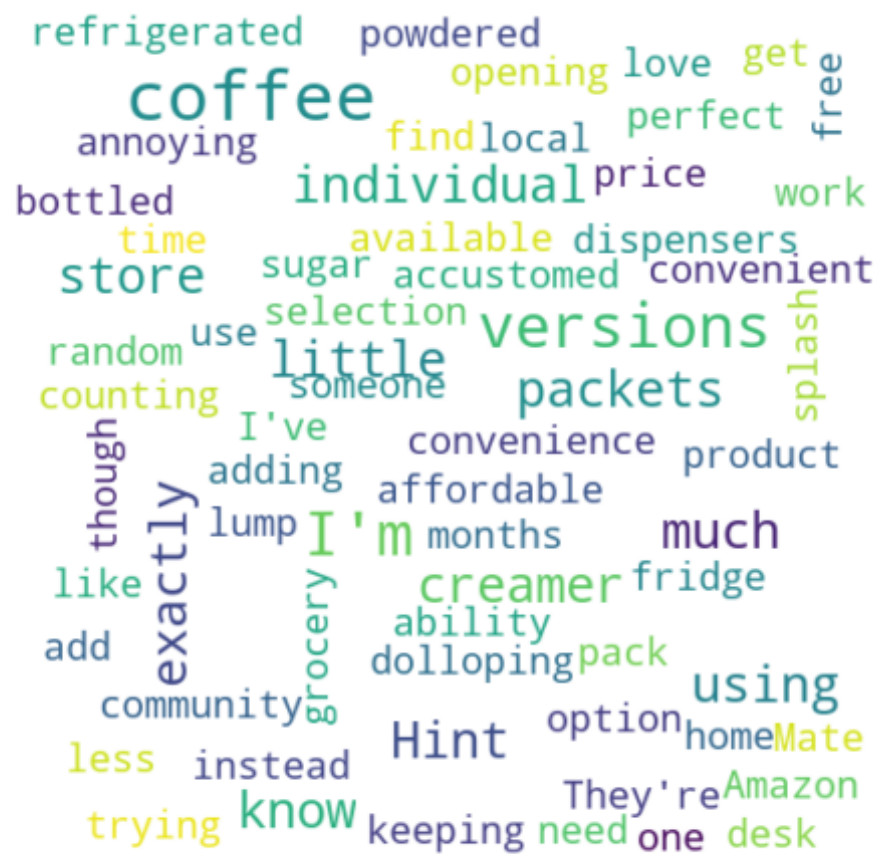
```
In [3]: # Please write all the code with proper documentation
```

```
In [39]: from wordcloud import WordCloud
print("cluster-1")
for i in range(2):
    print("review",i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10,max_font_size=40).generate(c1[i])
    # Display the generated image:
    plt.figure(figsize = (8, 8), facecolor = None)
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis("off")
    plt.show()
```

```
cluster-1
review 0
```







Observation: Cluster-1 is about grocery items like coffee powder.

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

```
cluster-2
review 0
```

laxative not  
tried  
tea consumed  
cleansing  
unique flavor  
daily I've  
one Good  
effective  
harsh needed  
regular system

review 1

future every  
trying



Observation: Cluster-2 is about the different tea flavor.

```
In [41]: from wordcloud import WordCloud
print("cluster-3")
for i in range(2):
    print("review", i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10, max_font_size=40).generate(c3[i])
    # Display the generated image:
    plt.figure(figsize = (8, 8), facecolor = None)
    plt.imshow(wordcloud, interpolation='bilinear')
```

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

cluster-3  
review 0

stinky VICTOR WWW  
available course  
traps M380 nearby  
Pretty BAIT  
product amazon  
right unreal FLY  
M502 REFILL  
com genocide  
http B00004RBDY  
total  
MAGNET 음

review 1

seasons  
fly  
Can't  
Great bait beat  
product  
used  
Victor

Observation:Cluster-3 is about review about a product may be victor

```
In [42]: from wordcloud import WordCloud
print("cluster-4")
for i in range(2):
    print("review",i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10,max_font_size=40).generate(c4[i])
    # Display the generated image:
    plt.figure(figsize = (8, 8), facecolor = None)
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis("off")
    plt.show()
```

cluster-4  
review 0

chicken find  
anymore  
bad buying wont take  
good going  
one product made  
imports loves  
isnt China dogs  
hard know  
till  
USA  
chances  
review 1  
dogs  
regarding safe  
pet



satisfied love  
saw  
attached China store  
tag made

Observation: Cluster-4 is about the pets food quality in different countries

### [5.1.5] Applying K-Means Clustering on AVG W2V, SET 3

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

```
In [43]: X = final["CleanedText"]  
print("shape of X:", X.shape)
```

shape of X: (46071,)

```
In [44]: # Train your own Word2Vec model using your own text corpus  
i=0  
list_of_sentence=[]  
for sentence in X:  
    list_of_sentence.append(sentence.split())
```

```
In [45]: 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_sentence,min_count=5,size=50, workers=4)
            print(w2v_model.wv.most_similar('great'))
            print('='*50)
            print(w2v_model.wv.most_similar('worst'))

        [('awesome', 0.8575524091720581), ('fantastic', 0.8252902030944824),
        ('good', 0.8243530988693237), ('wonderful', 0.7938828468322754), ('terr
        ific', 0.788963794708252), ('amazing', 0.7829962372779846), ('excellen
        t', 0.7786347270011902), ('perfect', 0.7528365254402161), ('nice', 0.69
        10187005996704), ('decent', 0.6888670325279236)]

        =====
        [('nastiest', 0.7416113018989563), ('greatest', 0.7266179919242859),
        ('best', 0.7094963788986206), ('tastiest', 0.6527155637741089), ('awfu
        l', 0.6511404514312744), ('horrible', 0.6335446834564209), ('experience
        d', 0.6327458620071411), ('closest', 0.6286952495574951), ('ive', 0.617
        1221733093262), ('eaten', 0.6110190749168396)]
```

```
In [46]: 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 12798
        sample words ['dogs', 'loves', 'chicken', 'product', 'china', 'wont',
        'buying', 'anymore', 'hard', 'find', 'products', 'made', 'usa', 'one',
        'isnt', 'bad', 'good', 'take', 'chances', 'till', 'know', 'going', 'imp
        orts', 'love', 'saw', 'pet', 'store', 'tag', 'attached', 'regarding',
        'satisfied', 'safe', 'available', 'victor', 'traps', 'unreal', 'cours
        e', 'total', 'fly', 'pretty', 'stinky', 'right', 'nearby', 'used', 'bai
        t', 'seasons', 'ca', 'not', 'beat', 'great']
```

```
In [47]: # average Word2Vec
        # compute average word2vec for each review.
        sent_vectors = []; # the avg-w2v for each sentence/review is stored in
```

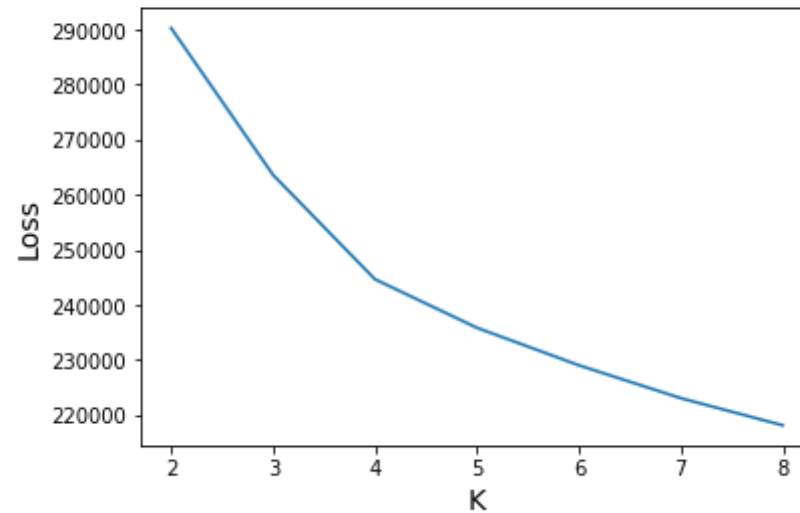
```
100%|███████████████████████████████████████████████████████████████████████████████  
██████████ | 46071/46071 [03:08<00:00, 244.72it/s]
```

46071  
50

```
In [35]: 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(sent_vectors)
    loss.append(k_means.inertia )
```

```
In [36]: plt.plot(k,loss)
plt.xlabel('K',size=14)
plt.ylabel('Loss',size=14)
plt.title('Loss VS K \n',size=18)
plt.show()
```

Loss VS K



```
In [48]: avgw2v_optimalk = 4
k_means = KMeans(n_clusters=avgw2v_optimalk, n_jobs=-1).fit(sent_vectors)
```

```
In [49]: reviews = final['Text'].values
c1 = []
c2 = []
c3 = []
c4 = []

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])
    elif k_means.labels_[i] == 2:
        c3.append(reviews[i])
    else:
```

```
c4.append(reviews[i])
```

### [5.1.6] Wordclouds of clusters obtained after applying k-means on AVG W2V SET 3

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

```
In [50]: from wordcloud import WordCloud
print("cluster-1")
for i in range(2):
    print("review",i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10,max_font_size=40).generate(c1[i])
    # Display the generated image:
    plt.figure(figsize = (8, 8), facecolor = None)
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis("off")
    plt.show()
```

```
cluster-1
review 0
```

unique Good  
system needed flavor  
I've regular daily tried  
tea  
laxative consumed  
not effective  
harsh  
one  
cleansing

review 1

new  
batches  
others  
strong  
need  
product  
careful  
stronger dosage

Observation:Cluster-1 is about the different tea flavor

```
In [51]: from wordcloud import WordCloud
print("cluster-2")
for i in range(2):
    print("review",i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10,max_font_size=40).generate(c2[i])
    # Display the generated image:
    plt.figure(figsize = (8, 8), facecolor = None)
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis("off")
    plt.show()
```

cluster-2  
review 0



not really  
bowl either pretty  
fish house quickly  
favorite cat  
want another  
clean one month get  
salmon bits  
outstanding  
Try definite make  
ground sliced times winner  
likes version still  
licks

review 1

market much  
ground better buy  
No tastes  
touch things  
leftovers bad  
even iety must  
K



Observation: Cluster-2 is about pet foods for dogs, cats, fish etc

```
In [52]: from wordcloud import WordCloud
print("cluster-3")
for i in range(2):
    print("review",i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10, max_font_size=40).generate(c3[i])
    # Display the generated image:
    plt.figure(figsize = (8, 8), facecolor = None)
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis("off")
    plt.show()
```

```
cluster-3
review 0
```

A word cloud visualization of text data. The words are arranged in a circular pattern, with their size and color indicating their frequency or importance. The word 'product' is the largest and most prominent, oriented vertically in the center. Other significant words include 'chicken', 'wont', 'China', 'product', 'hard', 'find', 'going', 'good', 'anymore', 'isnt', 'bad', 'till', 'USA', 'loves', 'made', 'buying', 'know', 'one', 'imports', 'dogs', 'chances', and 'review 1'. The colors range from green to yellow, with some words in purple.

USA till bad  
isnt  
take  
loves  
made buying chicken  
know  
one wont  
China  
imports dogs  
chances  
product  
hard  
find  
good  
going  
anymore

review 1



Observation: Cluster-3 is about the pets food quality in different countries

```
In [53]: from wordcloud import WordCloud
print("cluster-4")
for i in range(2):
    print("review", i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
```

```

stopwords = stopwords,
min_font_size = 10,max_font_size=40).generate(c4[i])
# Display the generated image:
plt.figure(figsize = (8, 8), facecolor = None)
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()

```

cluster-4  
review 0



review 1

Can't  
bait  
Victor  
used  
Great  
product  
seasons  
fly  
beat

Observation: Cluster-4 is about review about a product may be victor

#### [5.1.7] Applying K-Means Clustering on TFIDF W2V, SET 4

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

```
In [54]: X = final["CleanedText"]
```

```
print("shape of X:", X.shape)
```

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

```
In [55]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer()
tf_idf_matrix = model.fit_transform(X)
# 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 [56]: # 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]
            #
            # 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
    tfidf_sent_vectors.append(sent_vec)
    row += 1
print(len(tfidf_sent_vectors))
```

46071

```
In [41]: plt.plot(k,loss)
plt.xlabel('K',size=14)
plt.ylabel('Loss',size=14)
plt.title('Loss VS K \n',size=18)
plt.show()
```





```
_sent_vectors)
```

```
In [58]: reviews = final['Text'].values
c1 = []
c2 = []
c3 = []
c4 = []

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])
    elif k_means.labels_[i] == 2:
        c3.append(reviews[i])
    else:
        c4.append(reviews[i])
```

### [5.1.8] Wordclouds of clusters obtained after applying k-means on TFIDF W2V SET 4

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

```
In [59]: from wordcloud import WordCloud
print("cluster-1")
for i in range(2):
    print("review",i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10,max_font_size=40).generate(c1[i])
    # Display the generated image:
    plt.figure(figsize = (8, 8), facecolor = None)
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis("off")
    plt.show()
```

cluster-1  
review 0



A word cloud visualization for cluster-1 review 0. The words are arranged in a circular pattern, with 'product' and 'chicken' being the most prominent. Other visible words include 'anytime', 'USA', 'one', 'China', 'buying', 'going', 'wont', 'bad', 'take', 'made', 'hard', 'good', 'till', 'imports', 'know', 'dogs', 'chances', 'loves', 'find', 'isnt', and 'chicken'.

anytime  
USA  
one  
China  
buying  
going  
wont  
bad  
take  
made  
hard  
good  
till  
imports  
know  
dogs  
chances  
loves  
find  
isnt  
chicken

review 1



Observation: Cluster-1 is about the pets food quality in different countries

```
In [60]: from wordcloud import WordCloud
print("cluster-2")
for i in range(2):
    print("review",i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10,max_font_size=40).generate(c2[i])
```

```
# Display the generated image:  
plt.figure(figsize = (8, 8), facecolor = None)  
plt.imshow(wordcloud, interpolation='bilinear')  
plt.axis("off")  
plt.show()
```

cluster-2  
review 0



review 1

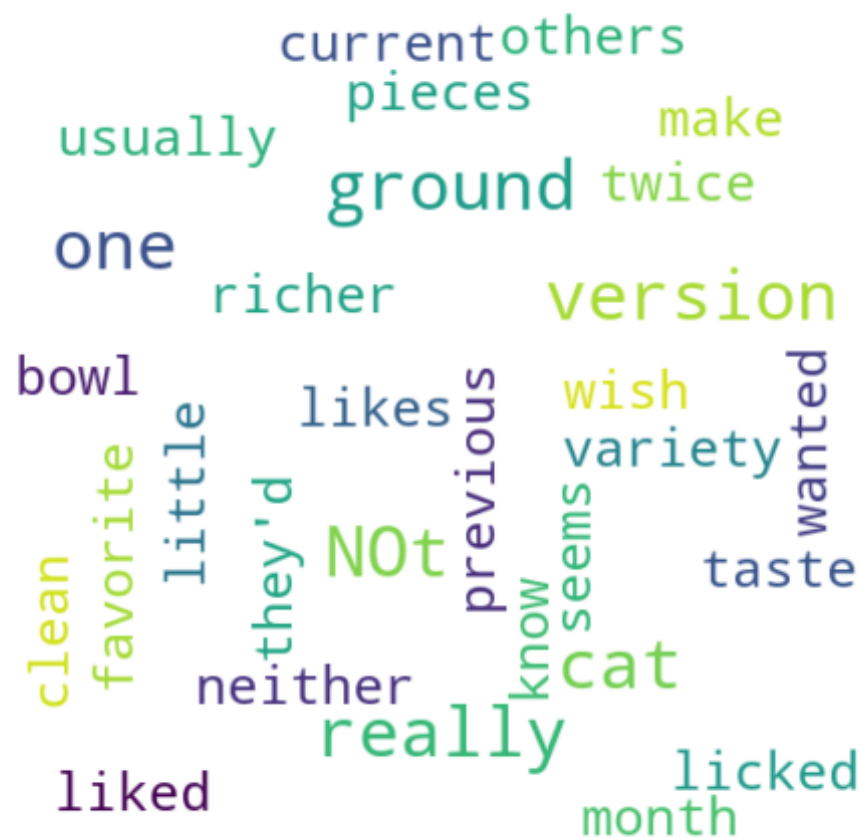


Observation: Cluster-2 is about review about a product may be victor

```
In [61]: from wordcloud import WordCloud
print("cluster-3")
for i in range(2):
    print("review", i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
```

cluster-3  
review 0



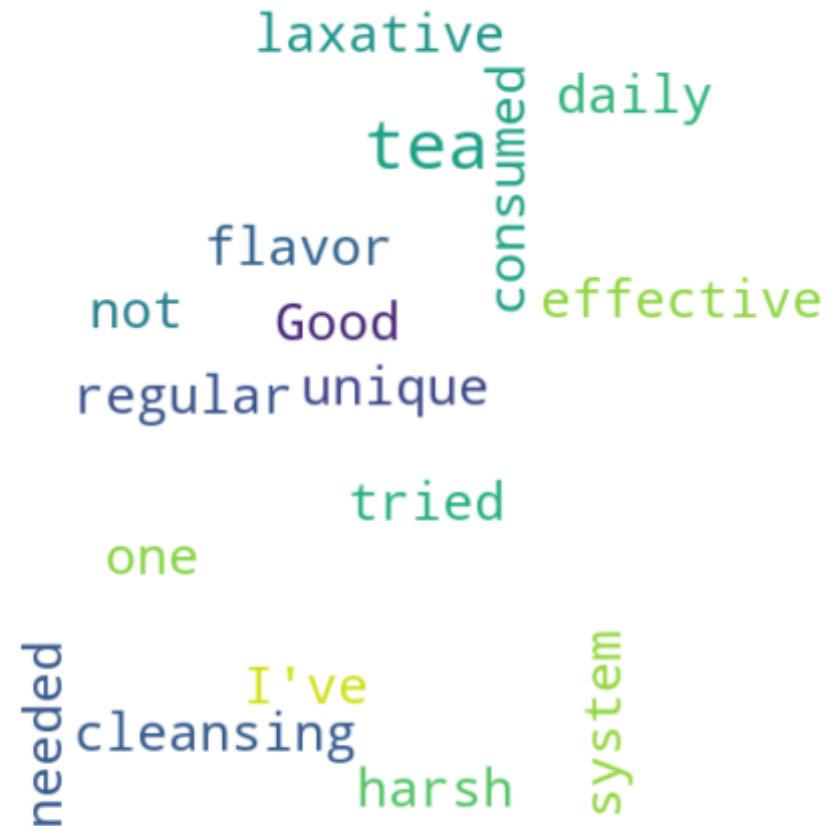


Observation: Cluster-3 is about pet foods for cats, fish etc

```
In [63]: from wordcloud import WordCloud
print("cluster-4")
for i in range(2):
    print("review", i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10, max_font_size=40).generate(c4[i])
```

```
# Display the generated image:  
plt.figure(figsize = (8, 8), facecolor = None)  
plt.imshow(wordcloud, interpolation='bilinear')  
plt.axis("off")  
plt.show()
```

cluster-4  
review 0



review 1





```
final = final.iloc[:5000,:]  
print(final.shape)  
(5000, 11)
```

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

shape of X: (5000,)

```
# Train your own Word2Vec model using your own text corpus
i=0
list_of_sentence=[]
for sentence in X:
    list_of_sentence.append(sentence.split())
```

```
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'))

[('think', 0.9926537871360779), ('pretty', 0.9907428026199341), ('good', 0.9905633926391602), ('texture', 0.9900360107421875), ('want', 0.9897626042366028), ('makes', 0.9894203543663025), ('well', 0.9891479015350342), ('stuff', 0.9887492656707764), ('something', 0.9887170195579529), ('nice', 0.9882326126098633)]

=====

[('ones', 0.9995308518409729), ('uses', 0.99951171875), ('none', 0.9994901418685913), ('pain', 0.9994831681251526), ('end', 0.9994723200798035), ('name', 0.9994722604751587), ('nearly', 0.9994550347328186), ('wife', 0.9994549999999999)]
```

```
e', 0.9994481801986694), ('friend', 0.9994431138038635), ('type', 0.999433159828186)]
```

```
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 4108
sample words ['dogs', 'loves', 'chicken', 'product', 'china', 'wont',
'buying', 'anymore', 'hard', 'find', 'products', 'made', 'usa', 'one',
'isnt', 'bad', 'good', 'take', 'chances', 'till', 'know', 'going', 'lov
e', 'saw', 'pet', 'store', 'tag', 'attached', 'regarding', 'satisfied',
'safe', 'available', 'course', 'total', 'fly', 'pretty', 'stinky', 'rig
ht', 'nearby', 'used', 'ca', 'not', 'beat', 'great', 'received', 'shippm
ent', 'could', 'hardly', 'wait', 'try']
```

```
In [30]: 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]))
         from sklearn.model_selection import train_test_split
```

```
100%|██████████| 5000/5000 [00:09<00:00, 517.96it/s]
```

5000

```
In [31]: from sklearn.cluster import AgglomerativeClustering
         agc = AgglomerativeClustering(n_clusters=4).fit(sent_vectors)
```

```
In [33]: reviews = final['Text'].values
         c1 = []
         c2 = []
         c3 = []
         c4 = []

         for i in range(agc.labels_.shape[0]):
             if agc.labels_[i] == 0:
                 c1.append(reviews[i])
             elif agc.labels_[i] == 1:
                 c2.append(reviews[i])
             elif agc.labels_[i] == 2:
                 c3.append(reviews[i])
             else :
                 c4.append(reviews[i])
```

### [5.2.2] Wordclouds of clusters obtained after applying Agglomerative Clustering on AVG W2V SET 3

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

```
In [34]: from wordcloud import WordCloud
         print("cluster-1")
         for i in range(2):
             print("review",i)
             wordcloud = WordCloud(width = 500, height = 500,
                                   background_color = 'white',
                                   stopwords = stopwords,
                                   min_font_size = 10,max_font_size=40).generate(c1[i])
```

```
# Display the generated image:  
plt.figure(figsize = (8, 8), facecolor = None)  
plt.imshow(wordcloud, interpolation='bilinear')  
plt.axis("off")  
plt.show()
```

cluster-1  
review 0



review 1

Can't product  
seasons Great  
beat  
fly Victor  
used bait

Observation: Cluster-1 is about food product quality in different countries

```
In [35]: from wordcloud import WordCloud
print("cluster-2")
for i in range(2):
    print("review", i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
```

```
min_font_size = 10,max_font_size=40).generate(c2[i])  
# Display the generated image:  
plt.figure(figsize = (8, 8), facecolor = None)  
plt.imshow(wordcloud, interpolation='bilinear')  
plt.axis("off")  
plt.show()
```

cluster-2  
review 0



review 1



Observation: Cluster-2 is about product victor

```
In [36]: from wordcloud import WordCloud
print("cluster-3")
for i in range(2):
    print("review", i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10, max_font_size=40).generate(c3[i])
```



```
# Display the generated image:  
plt.figure(figsize = (8, 8), facecolor = None)  
plt.imshow(wordcloud, interpolation='bilinear')  
plt.axis("off")  
plt.show()
```

cluster-3  
review 0

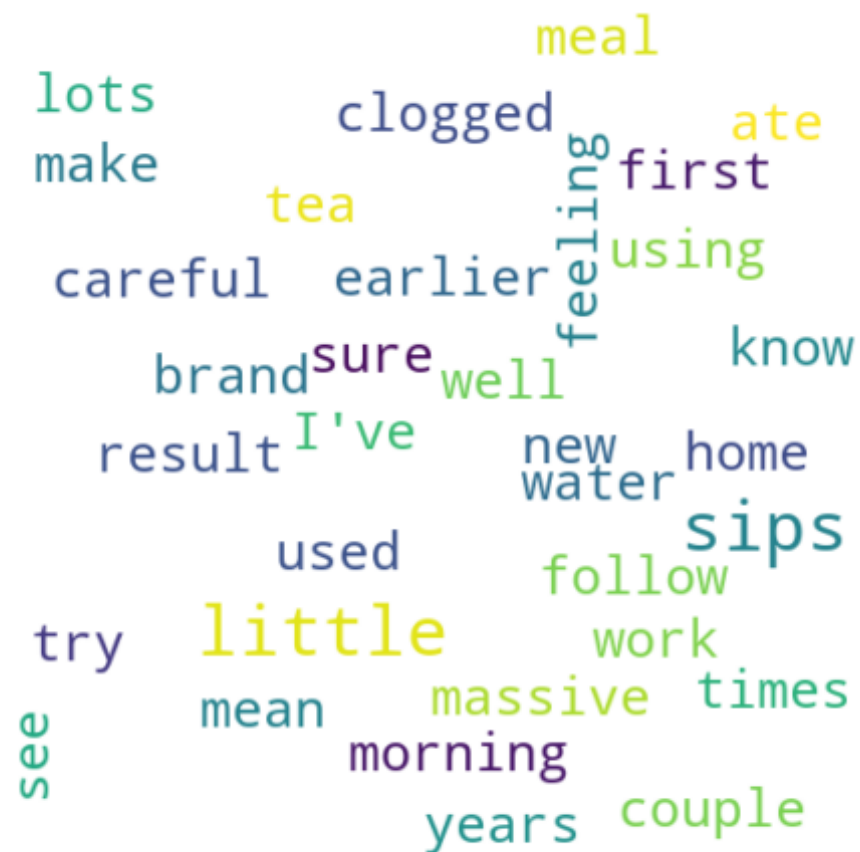


review 1



```
# Display the generated image:  
plt.figure(figsize = (8, 8), facecolor = None)  
plt.imshow(wordcloud, interpolation='bilinear')  
plt.axis("off")  
plt.show()
```

cluster-4  
review 0



review 1



Observation: Cluster-4 is about pet foods for fish, cat etc

### [5.2.3] Applying Agglomerative Clustering on TFIDF W2V, SET 4

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

```
In [38]: model = TfidfVectorizer()
```

```
tf_idf_matrix = model.fit_transform(X)
# 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 [39]: 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
    tfidf_sent_vectors.append(sent_vec)
    row += 1
print(len(tfidf_sent_vectors))
```

```
100%|██████████████████████████████████████████████████████████████████████████████|  
██████████ | 5000/5000 [00:50<00:00, 99.15it/s]
```

5000

```
In [40]: from sklearn.cluster import AgglomerativeClustering
         agc = AgglomerativeClustering(n_clusters=4).fit(tfidf_sent_vectors)
```

```
In [41]: reviews = final['Text'].values
c1 = []
c2 = []
c3 = []
c4 = []

for i in range(agc.labels_.shape[0]):
    if agc.labels_[i] == 0:
        c1.append(reviews[i])
    elif agc.labels_[i] == 1:
        c2.append(reviews[i])
    elif agc.labels_[i] == 2:
        c3.append(reviews[i])
    else :
        c4.append(reviews[i])
```

#### [5.2.4] Wordclouds of clusters obtained after applying Agglomerative Clustering on TFIDF W2V SET 4

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

```
In [42]: from wordcloud import WordCloud
print("cluster-1")
for i in range(2):
    print("review",i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10,max_font_size=40).generate(c1[i])
    # Display the generated image:
    plt.figure(figsize = (8, 8), facecolor = None)
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis("off")
    plt.show()
```

cluster-1  
review 0

China

made

store dogs

attached

safe

love

satisfied

tag

pet

regarding

saw

review 1



Observation: Cluster-1 is about pet product may be victor.

```
In [43]: from wordcloud import WordCloud
print("cluster-2")
for i in range(2):
    print("review",i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10,max_font_size=40).generate(c2[i])
```



```
# Display the generated image:  
plt.figure(figsize = (8, 8), facecolor = None)  
plt.imshow(wordcloud, interpolation='bilinear')  
plt.axis("off")  
plt.show()
```

cluster-2  
review 0



review 1

A word cloud visualization showing various terms. The words are arranged in a roughly circular pattern. The words include: 'dosage' (top left, dark blue), 'stronger' (top right, dark blue), 'new' (middle right, dark blue), 'product' (middle right, green), 'batches' (middle right, teal), 'strong' (bottom center, dark blue), 'careful' (bottom center, yellow-green), 'others' (middle left, teal), and 'need' (middle left, purple).

Observation: Cluster-2 is about food product quality in different countries.

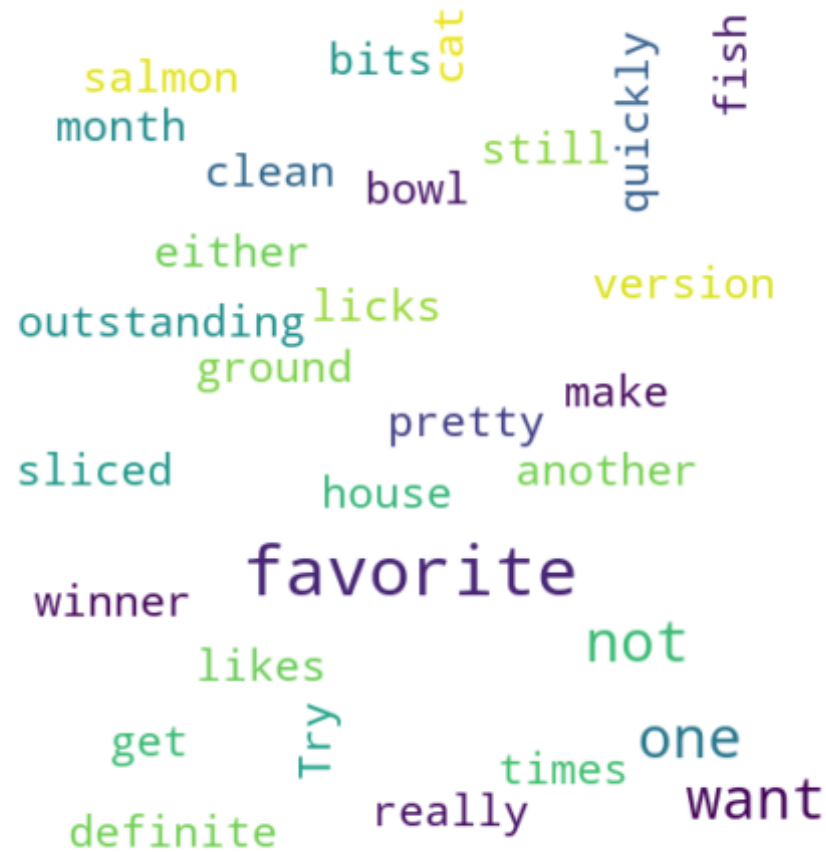
```
In [44]: from wordcloud import WordCloud
print("cluster-3")
for i in range(2):
    print("review", i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
```

```
min_font_size = 10,max_font_size=40).generate(c3[i])
# Display the generated image:
plt.figure(figsize = (8, 8), facecolor = None)
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```

cluster-3  
review 0



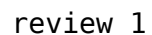
review 1

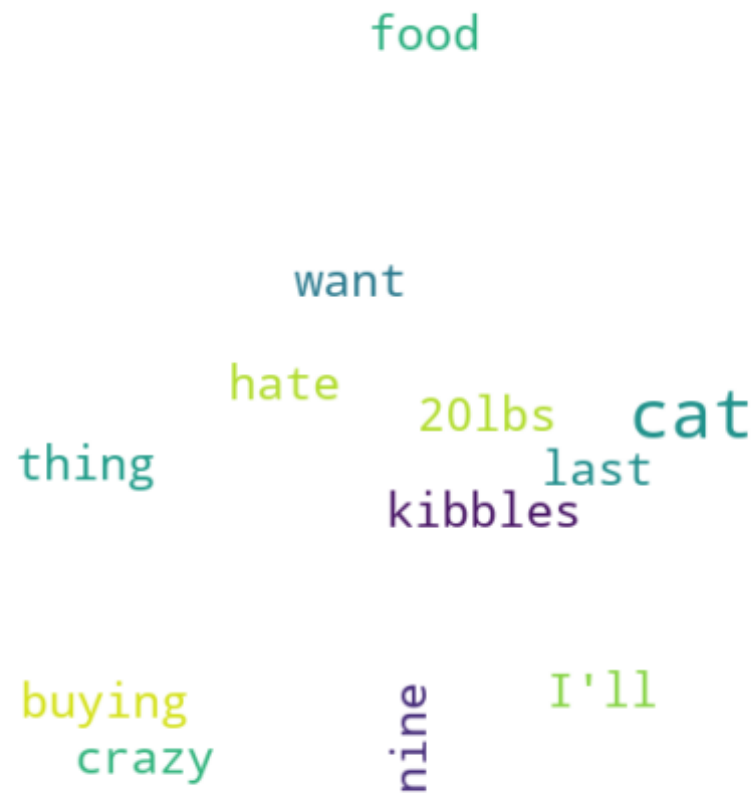


Observation: Cluster-3 is about pet foods for fish, cat etc.

```
In [45]: from wordcloud import WordCloud
print("cluster-4")
for i in range(2):
    print("review", i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10, max_font_size=40).generate(c4[i])
```

```
cluster-4
review 0
```





Observation: Cluster-4 is about protine quantity in food like 40lbs, 20lbs.

## [5.3] DBSCAN Clustering

### [5.3.1] Applying DBSCAN on AVG W2V, SET 3

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

```
In [69]: X = final["CleanedText"]
print("shape of X:", X.shape)
```

shape of X: (5000,)

```
In [70]: i=0
list_of_sentence=[]
for sentence in X:
    list_of_sentence.append(sentence.split())
```

```
In [71]: 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_sentence,min_count=5,size=50, workers=4)
    print(w2v_model.wv.most_similar('great'))
    print('='*50)
    print(w2v_model.wv.most_similar('worst'))
```

```
[('wonderful', 0.9898748993873596), ('nice', 0.989847719669342), ('make
s', 0.9887447953224182), ('bold', 0.988516628742218), ('yogi', 0.988253
116607666), ('spice', 0.9880969524383545), ('texture', 0.98800396919250
49), ('flavorful', 0.9879112243652344), ('smells', 0.9877873063087463),
('salty', 0.9877116084098816)]
```

=====

```
[('yeast', 0.9996388554573059), ('end', 0.9996042251586914), ('note',
0.9995812177658081), ('e', 0.9995694160461426), ('relieve', 0.999555945
3964233), ('type', 0.9995527267456055), ('truly', 0.9995508193969727),
('produced', 0.9995489716529846), ('variety', 0.9995481967926025), ('fr
iends', 0.9995433688163757)]
```

```
In [72]: 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 4108
sample words ['dogs', 'loves', 'chicken', 'product', 'china', 'wont',
'buying', 'anymore', 'hard', 'find', 'products', 'made', 'usa', 'one',
'isnt', 'bad', 'good', 'take', 'chances', 'till', 'know', 'going', 'love',
'saw', 'pet', 'store', 'tag', 'attached', 'regarding', 'satisfied',
'safe', 'available', 'course', 'total', 'fly', 'pretty', 'stinky', 'right',
'nearby', 'used', 'ca', 'not', 'beat', 'great', 'received', 'shipment',
'could', 'hardly', 'wait', 'try']
```

```
In [73]: 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, you
            # 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/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]))
        from sklearn.model_selection import train_test_split
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 5000/5000 [00:10<00:00, 474.72it/s]
```

```
5000
50
```

```
In [74]: sent_vector=np.array(sent_vectors)
```

```
In [75]: def n_neighbour(vectors , n):
        distance = []
```



```

for point in vectors:
    temp = np.sort(np.sum((vectors-point)**2,axis=1),axis=None)
    distance.append(temp[n])
return np.sqrt(np.array(distance))

```

```

In [76]: minpts = 2*50
dist = n_neighbour(sent_vector,minpts)
sorted_distance = np.sort(dist)
points = [i for i in range(sent_vector.shape[0])]

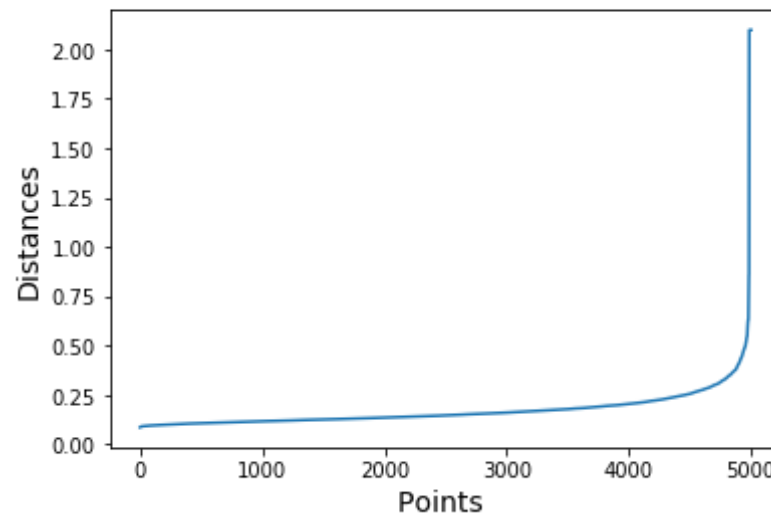
```

```

In [77]: plt.plot(points, sorted_distance)
plt.xlabel('Points ',size=14)
plt.ylabel('Distances',size=14)
plt.title('Distances VS Points Plot\n',size=18)
plt.show()

```

Distances VS Points Plot



```

In [78]: from sklearn.cluster import DBSCAN
dbs= DBSCAN(eps = 0.5, min_samples = minpts, n_jobs=-1).fit(sent_vector
)

```

```
In [79]: reviews = final['Text'].values
c=[]

for i in range(dbs.labels_.shape[0]):
    c.append(reviews[i])
```

### [5.3.2] Wordclouds of clusters obtained after applying DBSCAN on AVG W2V SET 3

```
In [2]: # Please write all the code with proper documentation
```

```
In [80]: from wordcloud import WordCloud
print("cluster")
for i in range(2):
    print("review",i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10,max_font_size=40).generate(c[i])
    # Display the generated image:
    plt.figure(figsize = (8, 8), facecolor = None)
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis("off")
    plt.show()
```

```
cluster
review 0
```



review 1

China  
saw dogs  
satisfied safe  
pet love  
tag store  
regarding made  
attached

Observation: Cluster is about pet food products quality in different countries.

### [5.3.3] Applying DBSCAN on TFIDF W2V, SET 4

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

```
In [81]: model = TfidfVectorizer()  
         tf_idf_matrix = model.fit_transform(X)
```

```
# 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 [82]: tfidf_feat = model.get_feature_names() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and ce
ll_val = tfidf

tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is st
ored in this list
row=0;
for sent in tqdm(list_of_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/r
view
    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
print(len(tfidf_sent_vectors))
```

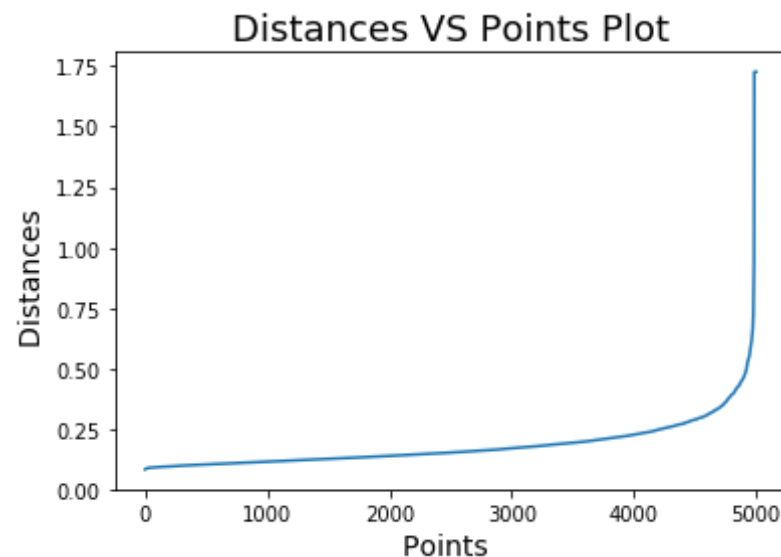
```
100% | ████████████████████████████████████████████████████████████  
████████████████████ | 5000/5000 [01:04<00:00, 77.38it/s]
```

5000

```
In [83]: tfidf_sent_vector=np.array(tfidf_sent_vectors)
```

```
In [84]: minpts = 2*50
dists = n_neighbour(tfidf_sent_vector,minpts)
sorted_distance = np.sort(dists)
point = [i for i in range(tfidf_sent_vector.shape[0])]
```

```
In [85]: plt.plot(points, sorted_distance)
plt.xlabel('Points ',size=14)
plt.ylabel('Distances ',size=14)
plt.title('Distances VS Points Plot',size=18)
plt.show()
```



```
In [86]: from sklearn.cluster import DBSCAN
dbs= DBSCAN(eps = 0.5, min_samples = minpts, n_jobs=-1).fit(tfidf_sent_vector)
```

```
In [87]: reviews = final['Text'].values
c = []
for i in range(dbs.labels_.shape[0]):
    c.append(reviews[i])
```

### [5.3.4] Wordclouds of clusters obtained after applying DBSCAN on TFIDF W2V SET 4

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

```
In [88]: from wordcloud import WordCloud
print("cluster")
for i in range(2):
    print("review",i)
    wordcloud = WordCloud(width = 500, height = 500,
                           background_color = 'white',
                           stopwords = stopwords,
                           min_font_size = 10,max_font_size=40).generate(c[i])
    # Display the generated image:
    plt.figure(figsize = (8, 8), facecolor = None)
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis("off")
    plt.show()
```

```
cluster
review 0
```

A word cloud featuring various words in different colors and sizes. The words are arranged in a somewhat circular pattern. The words include: find, bad, anymore, isnt, product, loves, made, one, China, USA, wont, know, hard, till, imports, chicken, take, going, buying, good, chances, dogs.

review 1



A word cloud visualization showing various terms related to pet food products quality. The words are arranged in a roughly circular pattern, with 'China' at the top, 'made' below it, 'dogs' to the right, 'saw' below 'dogs', 'safe' below 'saw', 'pet' to the left, 'store' to the right, 'attached' below 'safe', 'regarding' below 'attached', 'satisfied' below 'regarding', 'love' below 'satisfied', and 'tag' at the bottom. The words are in different colors: teal, yellow, green, blue, and purple.

China  
made  
dogs  
saw  
safe  
pet  
store  
attached  
regarding  
satisfied  
love  
tag

Observation: Cluster is about pet food products quality in different countries.

## [6] Conclusions

In [4]: *# Please compare all your models using Prettytable library.*

```
# You can have 3 tables, one each for kmeans, agglomerative and dbscan
```

```
In [89]: models = pd.DataFrame({'vectorizer': ['Kmeans with Bow', "Kmeans with TFIDF", "Kmeans with avg_w2v", "Kmeans with tfidf_avg_w2v"], 'Clustering': ["KMeans", "KMeans", "KMeans", "KMeans"], 'Number of cluster': [4, 4, 4, 4]}, columns = ["vectorizer", "Clustering", "Number of cluster"])
models
```

Out[89]:

	vectorizer	Clustering	Number of cluster
0	Kmeans with Bow	KMeans	4
1	Kmeans with TFIDF	KMeans	4
2	Kmeans with avg_w2v	KMeans	4
3	Kmeans with tfidf_avg_w2v	KMeans	4

```
In [92]: models = pd.DataFrame({'vectorizer': ["Agglomerative with avg_w2v", "Agglomerative with tfidf_avg_w2v"], 'Clustering': ["Hierarchical", "Hierarchical"], 'Number of cluster': [4, 4]}, columns = ["vectorizer", "Clustering", "Number of cluster"])
models
```

Out[92]:

	vectorizer	Clustering	Number of cluster
0	Agglomerative with avg_w2v	Hierarchical	4
1	Agglomerative with tfidf_avg_w2v	Hierarchical	4

```
In [93]: models = pd.DataFrame({'vectorizer': ["Dbscan with avg_w2v", "Dbscan with tfidf_avg_w2v"], 'Clustering': ["DBSCAN", "DBSCAN"], 'EPS': [.5, .5]}, columns = ["vectorizer", "Clustering", "EPS"])
models
```

Out[93]:

	vectorizer	Clustering	EPS
0	Dbscan with avg_w2v	DBSCAN	0.5

	vectorizer	Clustering	EPS
1	Dbscan with tfidf_avg_w2v	DBSCAN	0.5

In [ ]: