

Amazon Fine Food Reviews Analysis

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

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

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

Number of reviews: 568,454

Number of users: 256,059

Number of products: 74,258

Timespan: Oct 1999 - Oct 2012

Number of Attributes/Columns in data: 10

Attribute Information:

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

Objective:

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

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

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

▼ [1]. Reading Data

▼ [1.1] Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

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

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

```
1 %matplotlib inline
2 import warnings
3 warnings.filterwarnings("ignore")
4
5
6 import sqlite3
7 import pandas as pd
8 import numpy as np
9 import nltk
10 import string
11 import matplotlib.pyplot as plt
12 import seaborn as sns
13 from sklearn.feature_extraction.text import TfidfTransformer
14 from sklearn.feature_extraction.text import TfidfVectorizer
15
16 from sklearn.feature_extraction.text import CountVectorizer
17 from sklearn.metrics import confusion_matrix
18 from sklearn import metrics
19 from sklearn.metrics import roc_curve, auc
20 from nltk.stem.porter import PorterStemmer
21
22 import re
23 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
24 import string
25 from nltk.corpus import stopwords
26 from nltk.stem import PorterStemmer
27 from nltk.stem.wordnet import WordNetLemmatizer
28
29 from gensim.models import Word2Vec
30 from gensim.models import KeyedVectors
31 import pickle
32
33 from tqdm import tqdm
34 import os
```

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

👤 Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318

Enter your authorization code:
.....
Mounted at /content/drive

```
1 # using SQLite Table to read data.
2 con = sqlite3.connect('drive/My Drive/FFRDB/database.sqlite')
3 # filtering only positive and negative reviews i.e.
4 # not taking into consideration those reviews with Score=3
5 # SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000, will give top 500000 data points
6 # you can change the number to any other number based on your computing power
7
8 # filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3 LIMIT 50000
9 # for tsne assignment you can take 5k data points
10
11 filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3 LIMIT 100000"""
12
13 # Give reviews with Score>3 a positive rating(1), and reviews with a score<3 a negative ra
14 def partition(x):
15     if x < 3:
16         return 0
17     return 1
18
19 #changing reviews with score less than 3 to be positive and vice-versa
20 actualScore = filtered_data['Score']
21 positiveNegative = actualScore.map(partition)
22 filtered_data['Score'] = positiveNegative
23 print("Number of data points in our data", filtered_data.shape)
24 filtered_data.head(3)
```

👤 Number of data points in our data (100000, 10)

0	1	B001E4KFG0	A3SGXH7AUHU8GW		delmartian	1
1	2	B00813GRG4	A1D87F6ZCVE5NK		dll pa	0
2	3	B000LQOCHO	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"		1

▼ [2] Exploratory Data Analysis

```

1 #Sorting data according to ProductId in ascending order
2 sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True, inplace=False,
                                         na_position='last')

1 #Deduplication of entries
2 final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"}, keep='first')
3 final.shape

```

 (87775, 10)

```

1 final.sort_values('Time',inplace=True)
2 print(final.head(5))

```

	Id	...	Text
70688	76882	...	I bought a few of these after my apartment was...
1146	1245	...	This was a really good idea and the final prod...
1145	1244	...	I just received my shipment and could hardly w...
28086	30629	...	Nothing against the product, but it does bothe...
28087	30630	...	I love this stuff. It is sugar-free so it does...

[5 rows x 10 columns]

```

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

```

 87.775

```
1 final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```

1 #Before starting the next phase of preprocessing lets see the number of entries left
2 print(final.shape)
3
4 #How many positive and negative reviews are present in our dataset?
5 final['Score'].value_counts()

```

 (87773, 10)

1	73592
0	14181
Name: Score, dtype: int64	

▼ [3] Preprocessing

```
1 final=final.sample(50000)
```

```

1 # https://stackoverflow.com/a/47091490/4084039
2 import re

```

```

3 from bs4 import BeautifulSoup
4
5 def decontracted(phrase):
6     # specific
7     phrase = re.sub(r"won't", "will not", phrase)
8     phrase = re.sub(r"can't", "can not", phrase)
9
10    # general
11    phrase = re.sub(r"\n\t", " not", phrase)
12    phrase = re.sub(r"\ne", " are", phrase)
13    phrase = re.sub(r"\s", " is", phrase)
14    phrase = re.sub(r"\d", " would", phrase)
15    phrase = re.sub(r"\ll", " will", phrase)
16    phrase = re.sub(r"\t", " not", phrase)
17    phrase = re.sub(r"\ve", " have", phrase)
18    phrase = re.sub(r"\m", " am", phrase)
19
20    return phrase

```

```

1 # https://gist.github.com/sebleier/554280
2 # we are removing the words from the stop words list: 'no', 'nor', 'not'
3 # <br /><br /> ==> after the above steps, we are getting "br br"
4 # we are including them into stop words list
5 # instead of <br /> if we have <br/> these tags would have removed in the 1st step
6
7 stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ourselves',
8                 'you'll', "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'hi
9                 'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they
10                'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that"
11                'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had',
12                'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as',
13                'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through',
14                'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over',
15                'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any',
16                'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too',
17                's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now',
18                've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't",
19                'd', 'hadn', 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn',
20                'mustn', 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn
21                'won', "won't", 'wouldn', "wouldn't"])

```

```

1 # Combining all the above stundents
2 from tqdm import tqdm
3 preprocessed_reviews = []
4 # tqdm is for printing the status bar
5 for sentance in tqdm(final['Text'].values):
6     sentance = re.sub(r"http\S+", "", sentance)
7     sentance = BeautifulSoup(sentance, 'lxml').get_text()
8     sentance = decontracted(sentance)
9     sentance = re.sub("\S*\d\S*", "", sentance).strip()
10    sentance = re.sub('[^A-Za-z]+', ' ', sentance)

```

```

11 # https://gist.github.com/sebleier/554280
12 sentance = ' '.join(e.lower() for e in sentance.split() if e.lower() not in stopwords)
13 preprocessed_reviews.append(sentance.strip())

```

👤 100% |██████████| 50000/50000 [00:19<00:00, 2626.69it/s]

▼ [4] Featurization

▼ [4.1] BAG OF WORDS

```

1 count_vect = CountVectorizer(min_df=100)
2 BOW_50k = count_vect.fit_transform(preprocessed_reviews)
3 print("the type of count vectorizer ",type(BOW_50k))
4 print("the shape of out text BOW vectorizer ",BOW_50k.get_shape())
5 print("the number of unique words ", BOW_50k.get_shape()[1])

```

👤 the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer (50000, 2239)
the number of unique words 2239

▼ [4.3] TF-IDF

```

1 tf_idf_vect = TfidfVectorizer(min_df=100)
2 tfidf_50k=tf_idf_vect.fit_transform(preprocessed_reviews)
3
4 print("the type of count vectorizer ",type(tfidf_50k))
5 print("the shape of out text TFIDF vectorizer ",tfidf_50k.get_shape())
6 print("the number of unique words including both unigrams and bigrams ", tfidf_50k.get_sha

```

👤 the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer (50000, 2239)
the number of unique words including both unigrams and bigrams 2239

▼ [4.4] Word2Vec

```

1 # Train your own Word2Vec model using your own text corpus
2 i=0
3 list_of_sentance=[]
4 for sentance in preprocessed_reviews:
5     list_of_sentance.append(sentance.split())

```

```
1 # Using Google News Word2Vectors
2
3 # in this project we are using a pretrained model by google
4 # its 3.3G file, once you load this into your memory
5 # it occupies ~9Gb, so please do this step only if you have >12G of ram
6 # we will provide a pickle file which contains a dict ,
7 # and it contains all our corpus words as keys and model[word] as values
8 # To use this code-snippet, download "GoogleNews-vectors-negative300.bin"
9 # from https://drive.google.com/file/d/0B7XkCwpI5KDYNlNUTT1SS21pQmM/edit
10 # it's 1.9GB in size.
11
12
13 # http://kavita-ganesan.com/gensim-word2vec-tutorial-starter-code/#.W17SRFAzZPY
14 # you can comment this whole cell
15 # or change these variable according to your need
16
17 is_your_ram_gt_16g=False
18 want_to_use_google_w2v = False
19 want_to_train_w2v = True
20
21 if want_to_train_w2v:
22     # min_count = 5 considers only words that occurred at least 5 times
23     w2v_model=Word2Vec(list_of_sentence,min_count=5,size=50, workers=4)
24     print(w2v_model.wv.most_similar('great'))
25     print('='*50)
26     print(w2v_model.wv.most_similar('worst'))
27
28 elif want_to_use_google_w2v and is_your_ram_gt_16g:
29     if os.path.isfile('GoogleNews-vectors-negative300.bin'):
30         w2v_model=KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin',
31             print(w2v_model.wv.most_similar('great'))
32             print(w2v_model.wv.most_similar('worst'))
33     else:
34         print("you don't have Google's word2vec file, keep want_to_train_w2v = True, to tr
```



```
1 w2v_words = list(w2v_model.wv.vocab)
2 print("number of words that occurred minimum 5 times ",len(w2v_words))
3 print("sample words ", w2v_words[0:50])
```



number of words that occurred minimum 5 times 13478
sample words ['lipton', 'go', 'stix', 'iced', 'green', 'tea', 'mix', 'honey', 'mango',

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

▼ [4.4.1.1] Avg W2v

```

1 # average Word2Vec
2 # compute average word2vec for each review.
3 sent_vectors_50k = []; # the avg-w2v for each sentence/review is stored in this list
4 for sent in tqdm(list_of_sentance): # for each review/sentence
5     sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to cha
6     cnt_words =0; # num of words with a valid vector in the sentence/review
7     for word in sent: # for each word in a review/sentence
8         if word in w2v_words:
9             vec = w2v_model.wv[word]
10            sent_vec += vec
11            cnt_words += 1
12        if cnt_words != 0:
13            sent_vec /= cnt_words
14        sent_vectors_50k.append(sent_vec)
15 print(len(sent_vectors_50k))
16 print(len(sent_vectors_50k[0]))

```

 100% |██████████| 50000/50000 [01:29<00:00, 560.87it/s] 50000
50

▼ [4.4.1.2] TFIDF weighted W2v

```

1 # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
2 model = TfidfVectorizer()
3 tf_idf_matrix = model.fit_transform(preprocessed_reviews)
4 # we are converting a dictionary with word as a key, and the idf as a value
5 dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))

1 # TF-IDF weighted Word2Vec
2 tfidf_feat = model.get_feature_names() # tfidf words/col-names
3 # final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val = tfidf
4
5 tfidf_sent_vectors_50k = []; # the tfidf-w2v for each sentence/review is stored in this li
6 row=0;
7 for sent in tqdm(list_of_sentance): # for each review/sentence
8     sent_vec = np.zeros(50) # as word vectors are of zero length
9     weight_sum =0; # num of words with a valid vector in the sentence/review
10    for word in sent: # for each word in a review/sentence
11        if word in w2v_words and word in tfidf_feat:
12            vec = w2v_model.wv[word]
13            tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
14            # to reduce the computation we are
15            # dictionary[word] = idf value of word in whole courpus
16            # sent.count(word) = tf valeus of word in this review

```

```
17     tf_idf = dictionary[word]*(sent.count(word)/len(sent))
18     sent_vec += (vec * tf_idf)
19     weight_sum += tf_idf
20 if weight_sum != 0:
21     sent_vec /= weight_sum
22 tfidf_sent_vectors_50k.append(sent_vec)
23 row += 1
```

👤 100% |██████████| 50000/50000 [18:05<00:00, 46.08it/s]

[5] Assignment 10: K-Means, Agglomerative & DBSCAN Clustering

1. Apply K-means Clustering on these feature sets:

- **SET 1:**Review text, preprocessed one converted into vectors using (BOW)
- **SET 2:**Review text, preprocessed one converted into vectors using (TFIDF)
- **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

```
1 ## using upto 10 cluster centers
2 NumCenters = [2,3,4,5,6,7,8,9,10]
3
4 #defining a function to return inertia for every set of cluster center used
5
6 def Inertia(NumCenters,data):
7     from sklearn.cluster import KMeans
8     Inertia=[]
9     for i in NumCenters:
10         clf = KMeans(n_clusters=i, init='k-means++', verbose=5, n_jobs=-1)
11         clf.fit(data)
12         kmeans=clf.inertia_
13         Inertia.append(kmeans)
14     return Inertia
15
16
17
18 # plotting The Inertia vs K graph
19
20 def PlotElbow(NumCenters,Inertia):
21     plt.plot(NumCenters, Inertia)
22     plt.xlabel('K-values',size=15)
23     plt.ylabel('Inertia',size=15)
24     plt.title('Inertia VS K-values Plot\n',size=20)
25     plt.grid()
26     plt.show()
27
28 '''
29 Algorithm to find elbow of a graph is taken from the following questionare on Stackoverflo
30 #####
31 https://stackoverflow.com/questions/2018178/finding-the-best-trade-off-point-on-a-curve
32 #####
33 '''
34
35 # finding the Elbow of graph
36
37 def ElbowFinder(Inertia):
38     import numpy as np
39     import numpy.matlib
40     nPoints = len(Inertia)
41     allCoord = np.vstack((range(nPoints), Inertia)).T
42     np.array([range(nPoints), Inertia])
43     firstPoint = allCoord[0]
44     lineVec = allCoord[-1] - allCoord[0]
45     lineVecNorm = lineVec / np.sqrt(np.sum(lineVec**2))
```

```
46 vecFromFirst = allCoord - firstPoint
47 scalarProduct = np.sum(vecFromFirst * np.matlib.repmat(lineVecNorm, nPoints, 1), axis=1)
48 vecFromFirstParallel = np.outer(scalarProduct, lineVecNorm)
49 vecToLine = vecFromFirst - vecFromFirstParallel
50 distToLine = np.sqrt(np.sum(vecToLine ** 2, axis=1))
51 return np.argmax(distToLine)
52
53
54
55
56
57
58 from wordcloud import WordCloud, STOPWORDS
59 import matplotlib.pyplot as plt
60 stopwords = set(STOPWORDS)
61
62 # Generating word cloud function for a given dataset of str
63
64 def ShowWordcloud(data, title = None):
65     wordcloud = WordCloud(
66         background_color='white',
67         stopwords=stopwords,
68         max_words=200,
69         max_font_size=40,
70         scale=3,
71         random_state=1 # chosen at random by flipping a coin; it was heads
72     ).generate(str(data))
73
74     fig = plt.figure(1, figsize=(12, 12))
75     plt.axis('off')
76
77
78     if title:
79         fig.suptitle(title, fontsize=20)
80         fig.subplots_adjust(top=2.3)
81
82     plt.imshow(wordcloud)
83     plt.show()
84
85
86
87 # overall function to generate graph and word cloud
88
89 def PrintWordcloud(preprocessed_reviews, vectorizer):
90
91 # finding inertia and plotting it against different no of cluster
92     inertia=Inertia(NumCenters,vectorizer)
93     print('Plotting elbow graph')
94     PlotElbow(NumCenters,inertia)
95     print('\n\n')
96 # defining kmeans with best K
97     from sklearn.cluster import KMeans
```

```
98 clf= KMeans(n_clusters=ElbowFinder(inertia), n_jobs=-1)
99 model=clf.fit(vectorizer)
100 labels=model.labels_
101 # printing all clusters
102 unq_labels=np.unique(labels)
103 print('All clusters are',unq_labels)
104 print("\n")
105
106 corpus = preprocessed_reviews
107 # list of list containing reviews belonging to each cluster
108 ListOfReviewsClusters=[[[] for i in unq_labels]]
109 # segregating reviews from each cluster
110 for index,word in enumerate(corpus):
111     real_label = labels[index]
112     ListOfReviewsClusters[real_label].append(word)
113 # printing each review cloud
114 for i in range(len(list(unq_labels))):
115     print('word cloud for cluster', i+1)
116     ShowWordcloud(ListOfReviewsClusters[i])
117     print('\n\n')
118
```

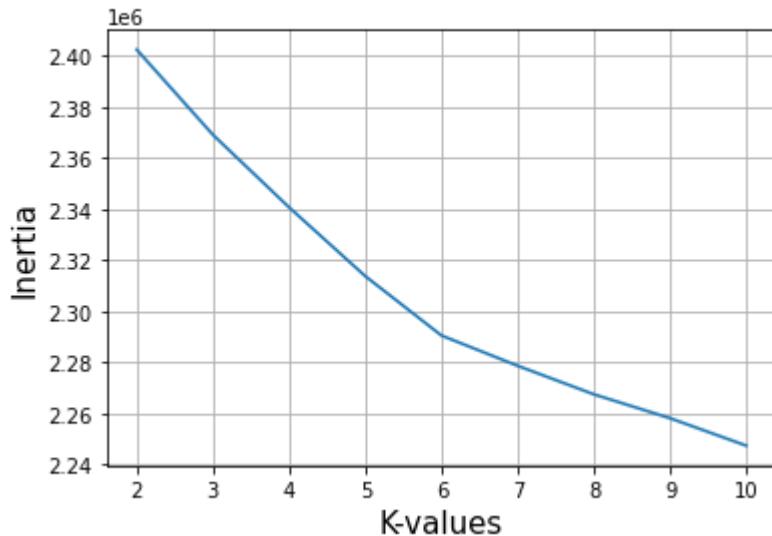
▼ [5.1.1] Applying K-Means Clustering on BOW, SET 1

1 PrintWordcloud(preprocessed_reviews,BOW_50k)



Plotting elbow graph

Inertia VS K-values Plot



All clusters are [0 1 2 3]

word cloud for cluster 1

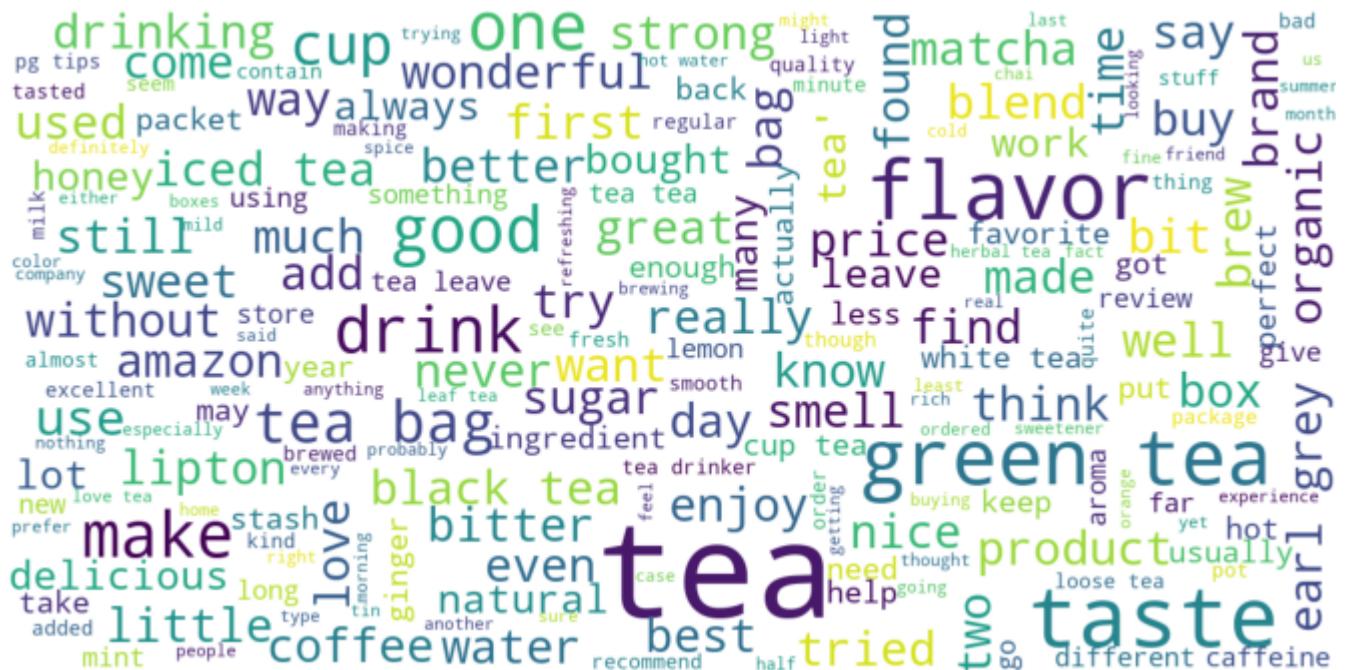


word cloud for cluster 2





word cloud for cluster 3



1. word cloud 1 represents good words which constitute good reviews.
 2. word cloud 2 represents mainly about coffee and its comparisions.
 3. word cloud 3 represents mainly about tea and its comparisions.
 4. word cloud 4 seems to cluster mainly noise.

Kind better look two

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

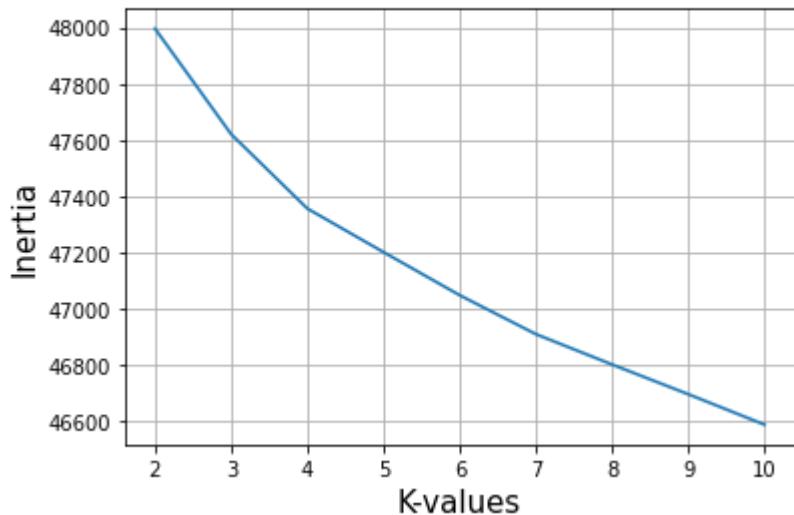
it's goal help right add

```
1 PrintWordcloud(preprocessed_reviews,tfidf_50k)
```



Plotting elbow graph

Inertia VS K-values Plot

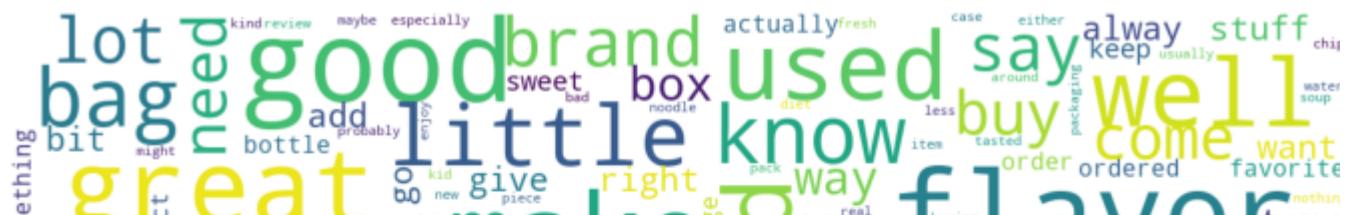


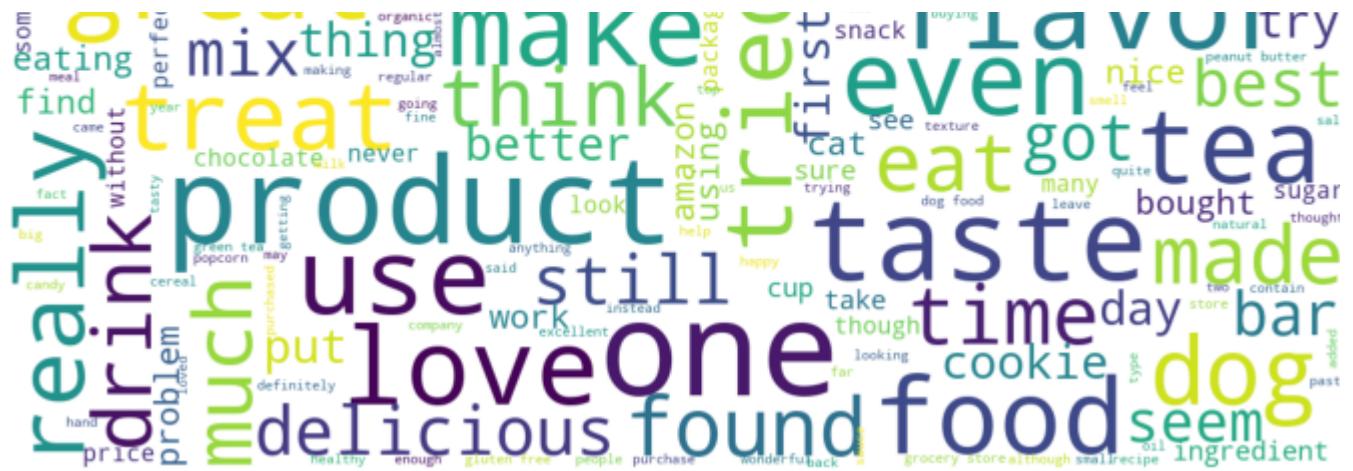
All clusters are [0 1]

word cloud for cluster 1



word cloud for cluster 2





1. word cloud 1 represents good review about Coffee.
 2. word cloud 2 represent in-general positive reviews.

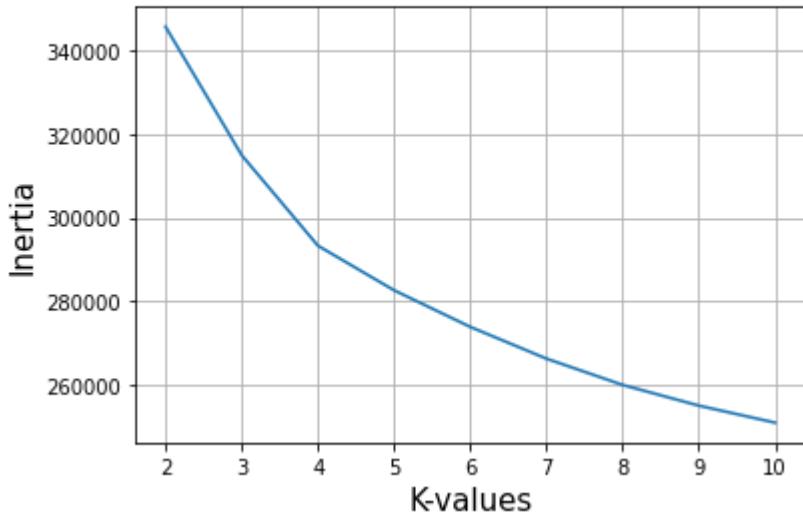
▼ [5.1.5] Applying K-Means Clustering on AVG W2V, SET 3

```
1 PrintWordcloud(preprocessed_reviews,sent_vectors_50k)
```



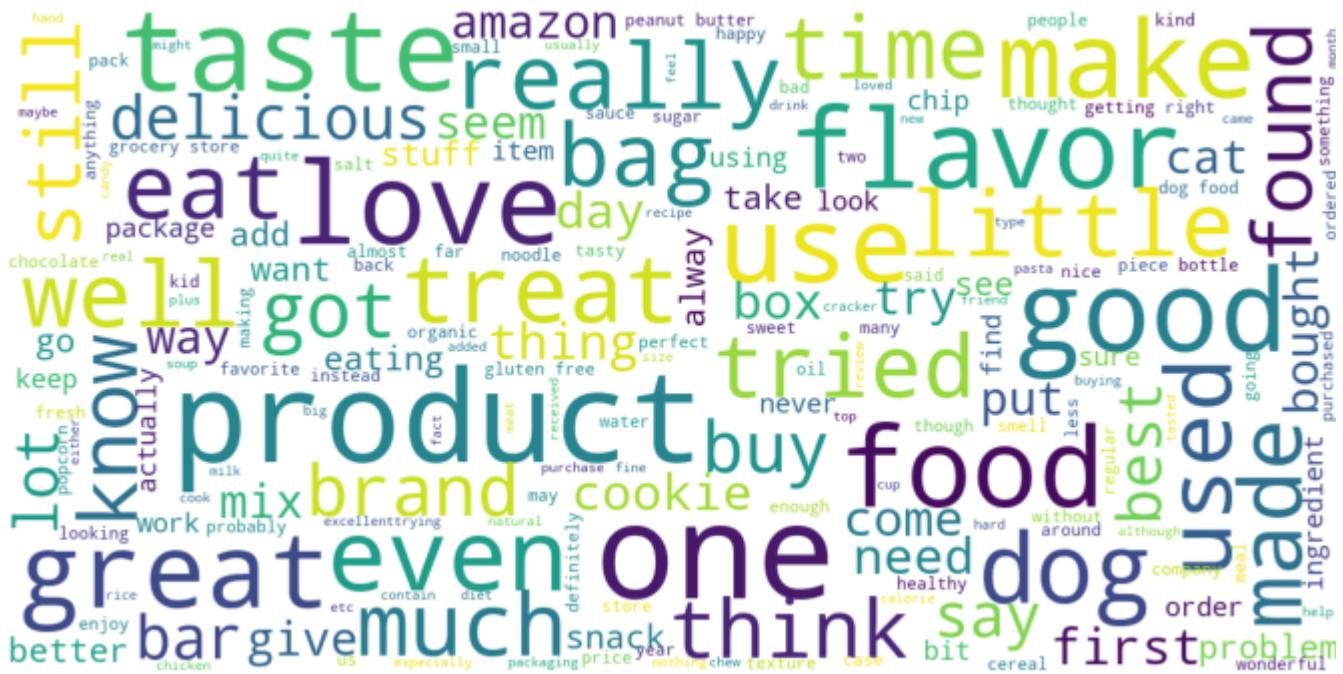
Plotting elbow graph

Inertia VS K-values Plot



All clusters are [0 1]

word cloud for cluster 1



1. word cloud 1 represent in-general positive reviews.
 2. word cloud 2 represents good review about Coffee and Tea.



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

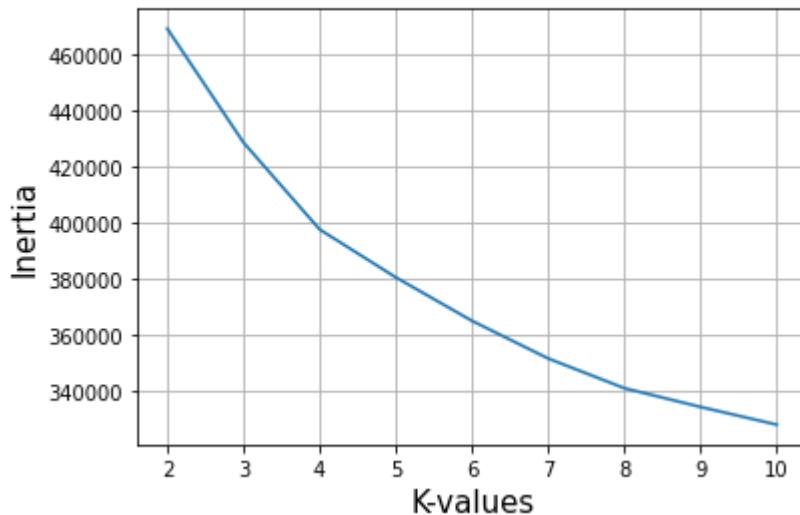


```
1 PrintWordcloud(preprocessed_reviews,tfidf_sent_vectors_50k)
```



Plotting elbow graph

Inertia VS K-values Plot



same as AVG W2V

1. word cloud 1 represents in general good reviews.
 2. word cloud 2 represents good review about Coffee and Tea.

twoOneTwo flavor for all your treats need cat

▼ [5.2] Agglomerative Clustering

size insight added

▼ [5.2.1] Applying Agglomerative Clustering on AVG W2V, SET 3

Instead of            

```
1 sent_vect=sent_vectors_50k[:100]
```

2 ppr=preprocessed reviews[:10000]

```
1 from wordcloud import WordCloud, STOPWORDS
```

```
3 import matplotlib.pyplot as plt
```

```
3 stopwords = set(STOPWORDS)
```

4

5 # Generating word cloud function for a given dataset of strings

5

```
3 def ShowWordCloud(data, title = None):
```

7

background color habitat

background_color = white
background_stayandahead

stopwords=stopwords, n_grams=300

max_words=200,
max_font_size=40

scale_3

14 random

<https://colab.research.google.com/drive/1ENc9YFcv05MIXyy-ZRPPBYMi71Of6MMRw/#scrollTo=HviwlwaeEsW&printMode=true>

```
15     ).generate(str(data))
16
17     fig = plt.figure(1, figsize=(12, 12))
18     plt.axis('off')
19
20
21     if title:
22         fig.suptitle(title, fontsize=20)
23         fig.subplots_adjust(top=2.3)
24
25     plt.imshow(wordcloud)
26     plt.show()
27
28
29
30 # overall function to generate graph and word cloud
31
32 def PrintWordcloud(preprocessed_reviews, vectorizer, k):
33
34
35     from sklearn.cluster import AgglomerativeClustering
36     clf = AgglomerativeClustering(n_clusters=k)
37     model=clf.fit(vectorizer)
38     labels=model.labels_
39
40     unq_labels=np.unique(labels)
41     # print('All clusters are',unq_labels)
42     # print("\n")
43
44     corpus = preprocessed_reviews
45     ListOfReviewsClusters=[[] for i in unq_labels]
46
47     for index,word in enumerate(corpus):
48         real_label = labels[index]
49         ListOfReviewsClusters[real_label].append(word)
50
51     for i in range(len(list(unq_labels))):
52         print('word cloud for cluster', i+1)
53         ShowWordcloud(ListOfReviewsClusters[i])
54         print('\n\n')
55
```

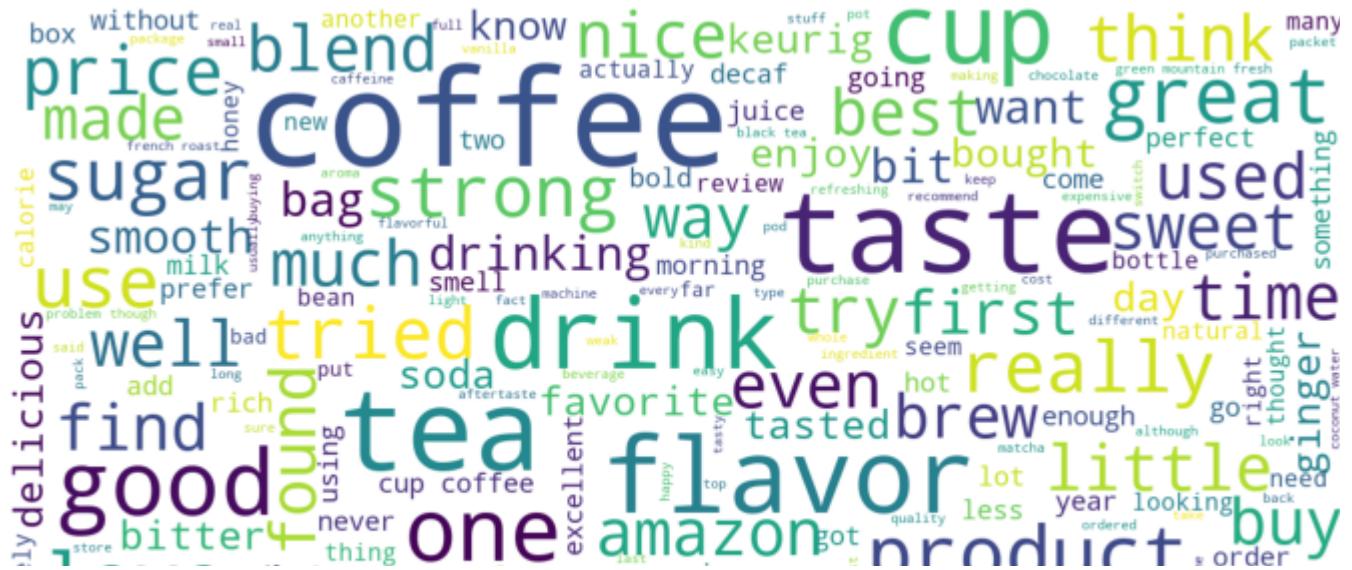
```
1 PrintWordcloud(ppr,sent_vect,2)
```



word cloud for cluster 1



word cloud for cluster 2



1. word cloud 1 represents general good reviews.
 2. word cloud 1 represents good review about Coffee and tea.

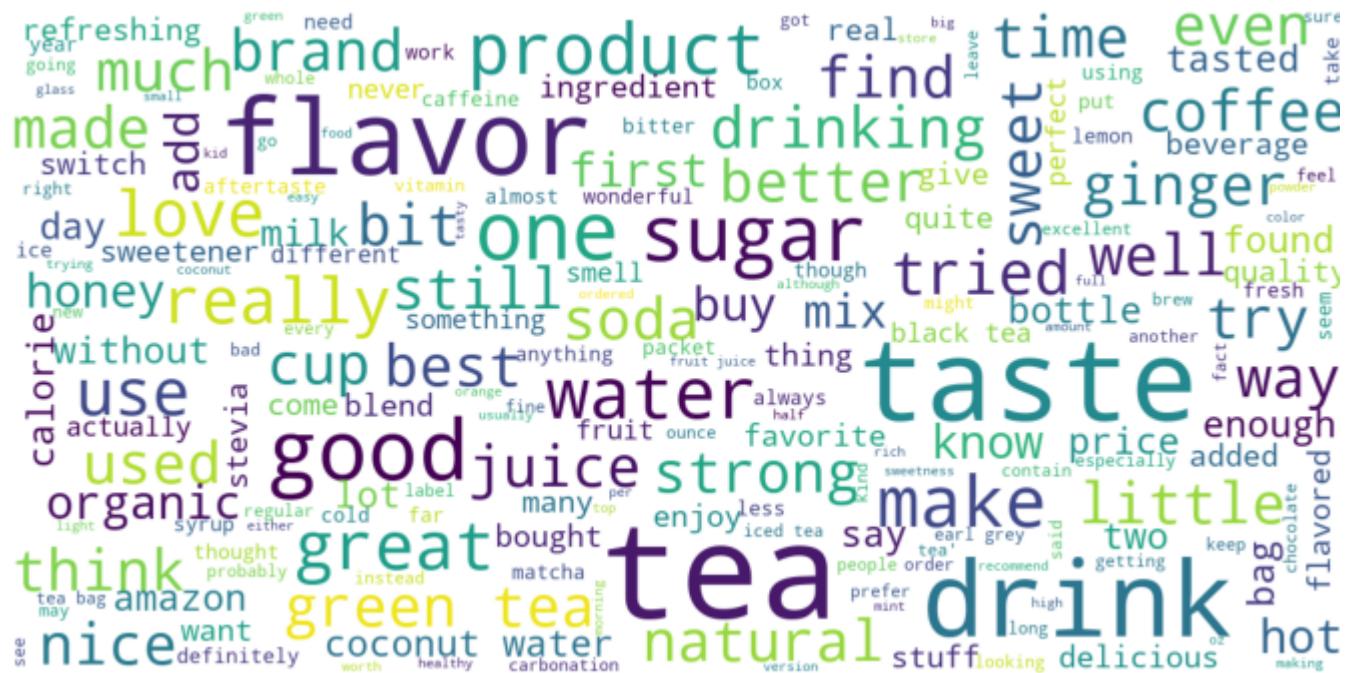
```
1 PrintWordcloud(ppr,sent_vect,5)
```



word cloud for cluster 1



word cloud for cluster 2



word cloud for cluster 3





word cloud for cluster 4



word cloud for cluster 5





1. word cloud 1 represents good review about the taste of products.
2. word cloud 2 represents good review about tea.
3. word cloud 3 represents good review about treats and dog/cat food.
4. word cloud 4 represents good review in general.
5. word cloud 5 represents good review about Coffee.

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

```
1 tfidf_sent_vect=tfidf_sent_vectors_50k[:10000]
```

```
1 PrintWordcloud(ppr,tfidf_sent_vect,2)
```



word cloud for cluster 1



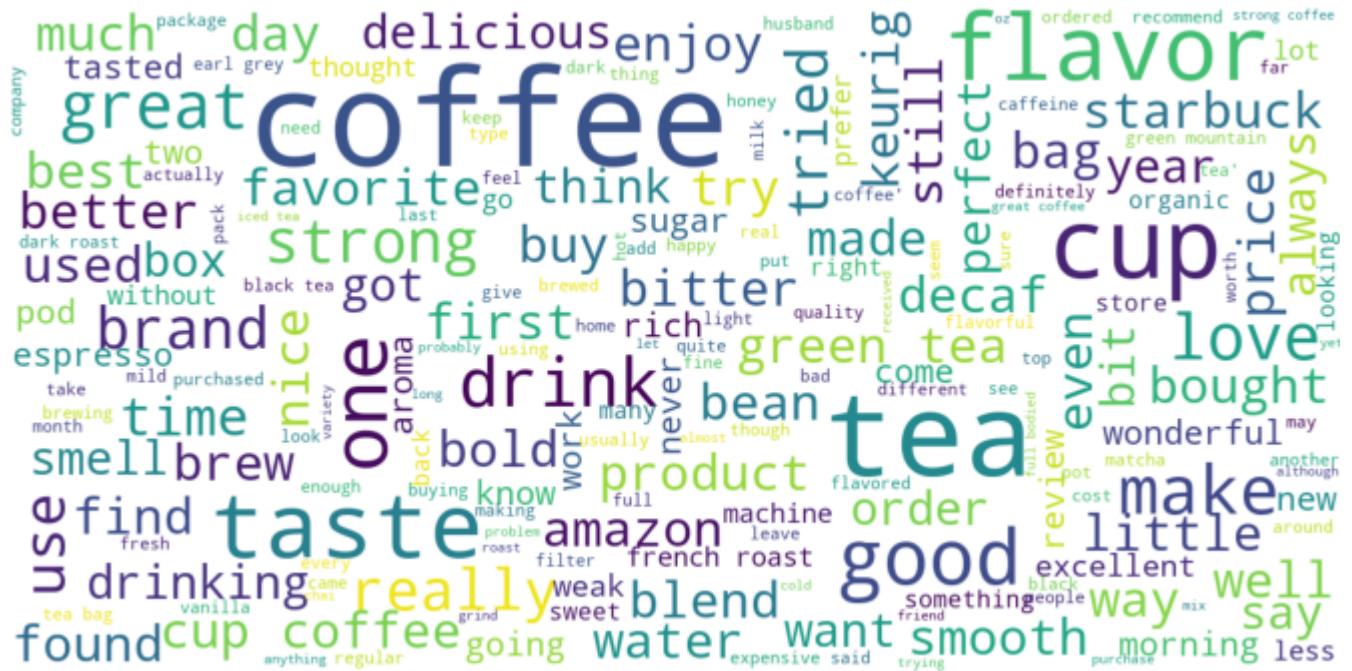
1. word cloud 1 represents good review ingeneral/Noise.
 2. word cloud 1 represents good review about Coffee and tea.

bit kid                              <img alt="leaf icon" data

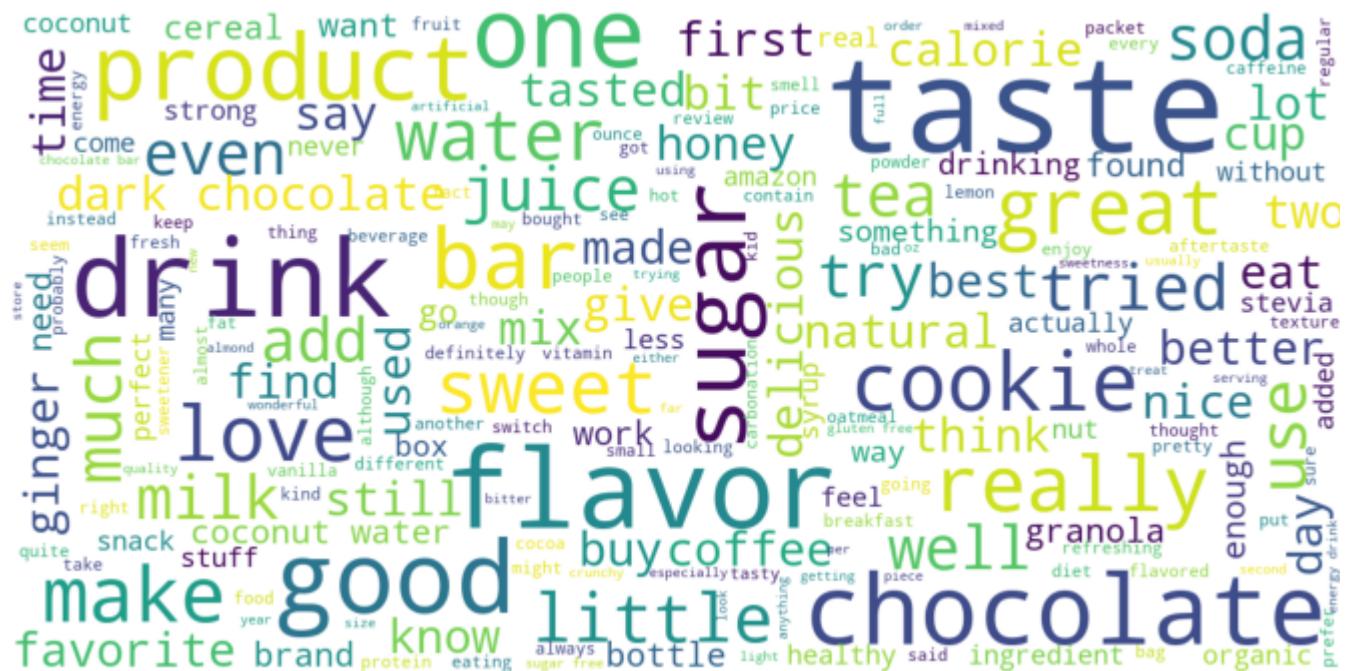
```
1 PrintWordcloud(ppr,tfidf_sent_vect,5)
```



word cloud for cluster 1



word cloud for cluster 2



word cloud for cluster 3





word cloud for cluster 4



1. word cloud 1 represents good review about Coffee and tea.
 2. word cloud 2 represents good review about the taste of products.
 3. word cloud 3 represents good review about treats and dog/cat food.
 4. word cloud 4 represents good review in general/Noise1.
 5. word cloud 5 represents good review in general/Noise2.

▼ [5.3] DBSCAN Clustering

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

```
1 ## using upto 10 cluster centers
2 NumCenters = [2,5]
3
4 #defining a function to return inertia for every set of cluster center used
5 !!!
```

```
7 """
8 def NeighborsDist(data,Min_pt=10):
9     from sklearn.neighbors import NearestNeighbors
10    nbrs = NearestNeighbors(n_neighbors=Min_pt).fit(data)
11    distances, indices = nbrs.kneighbors(data)
12    distances = np.sort(distances, axis=0)
13    distances = distances[:,1][::-1]
14    return distances
15
16
17 # plotting The Inertia vs K graph
18
19 def PlotElbow(data):
20     plt.plot(data)
21     plt.xlabel('Num_pts',size=15)
22     plt.ylabel('Distance',size=15)
23     plt.title('distance VS Num_pts Plot\n',size=20)
24     plt.grid()
25     plt.show()
26
27 """
28 Algorithm to find elbow of a graph is taken from the following questionare on Stackoverflo
29 #####
30 https://stackoverflow.com/questions/2018178/finding-the-best-trade-off-point-on-a-curve
31 #####
32 """
33
34 # finding the Elbow of graph
35
36 def ElbowFinder(Inertia):
37     import numpy as np
38     import numpy.matlib
39     nPoints = len(Inertia)
40     allCoord = np.vstack((range(nPoints), Inertia)).T
41     np.array([range(nPoints), Inertia])
42     firstPoint = allCoord[0]
43     lineVec = allCoord[-1] - allCoord[0]
44     lineVecNorm = lineVec / np.sqrt(np.sum(lineVec**2))
45     vecFromFirst = allCoord - firstPoint
46     scalarProduct = np.sum(vecFromFirst * np.matlib.repmat(lineVecNorm, nPoints, 1), axis=1)
47     vecFromFirstParallel = np.outer(scalarProduct, lineVecNorm)
48     vecToLine = vecFromFirst - vecFromFirstParallel
49     distToLine = np.sqrt(np.sum(vecToLine ** 2, axis=1))
50     return Inertia[np.argmax(distToLine)]
51
52
53
54
55
56
57 from wordcloud import WordCloud, STOPWORDS
```

```
58 import matplotlib.pyplot as plt
59 stopwords = set(STOPWORDS)
60
61 # Generating word cloud function for a given dataset of str
62
63 def ShowWordcloud(data, title = None):
64     wordcloud = WordCloud(
65         background_color='white',
66         stopwords=stopwords,
67         max_words=200,
68         max_font_size=40,
69         scale=3,
70         random_state=1 # chosen at random by flipping a coin; it was heads
71     ).generate(str(data))
72
73     fig = plt.figure(1, figsize=(12, 12))
74     plt.axis('off')
75
76
77     if title:
78         fig.suptitle(title, fontsize=20)
79         fig.subplots_adjust(top=2.3)
80
81     plt.imshow(wordcloud)
82     plt.show()
83
84
85
86 # overall function to generate graph and word cloud
87
88 def PrintWordcloud(preprocessed_reviews, vectorizer, samples=50):
89
90     nbr=NeighborsDist(vectorizer,Min_pt=samples)
91     print('Plotting Epsilon graph')
92     PlotElbow(nbr)
93     print('\n\n')
94     print('eps =',ElbowFinder(nbr))
95     print('\n\n')
96
97
98     from sklearn.cluster import DBSCAN
99
100    clf = DBSCAN(min_samples=samples,eps=ElbowFinder(nbr))
101    model=clf.fit(vectorizer)
102
103    labels=model.labels_
104
105    unq_labels=np.unique(labels)
106    print('All clusters are',unq_labels)
107    print("\n")
108
109    corpus = preprocessed_reviews
```

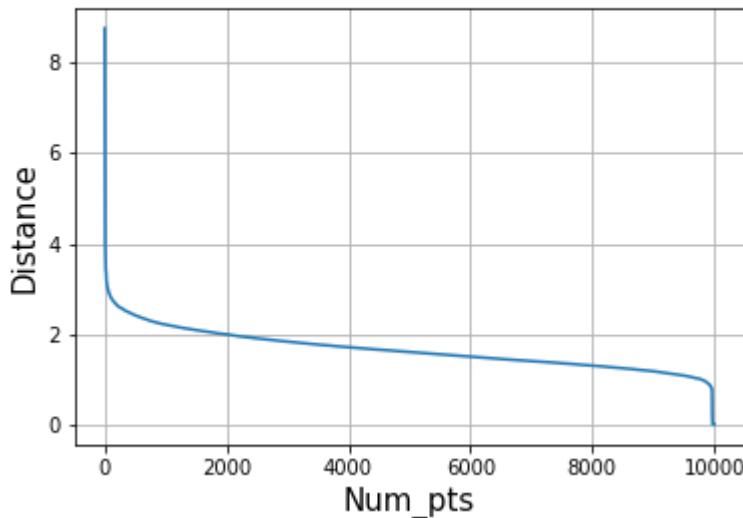
```
110 ListOfReviewsClusters=[[] for i in unq_labels]
111
112 for index,word in enumerate(corpus):
113     real_label = labels[index]
114     ListOfReviewsClusters[real_label].append(word)
115
116 for i in range(len(list(unq_labels))):
117     print('word cloud for cluster',i-1)
118     ShowWordcloud(ListOfReviewsClusters[i])
119     print('\n\n')
120
```

```
1 PrintWordcloud(ppr,sent_vect)
```



Plotting Epsilon graph

distance VS Num_pts Plot



eps = 2.5941881277686805

All clusters are [-1 0]

word cloud for cluster -1



1. word cloud -1 represents general Noise.
 2. word cloud 2 represents good review about coffee and tea.



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

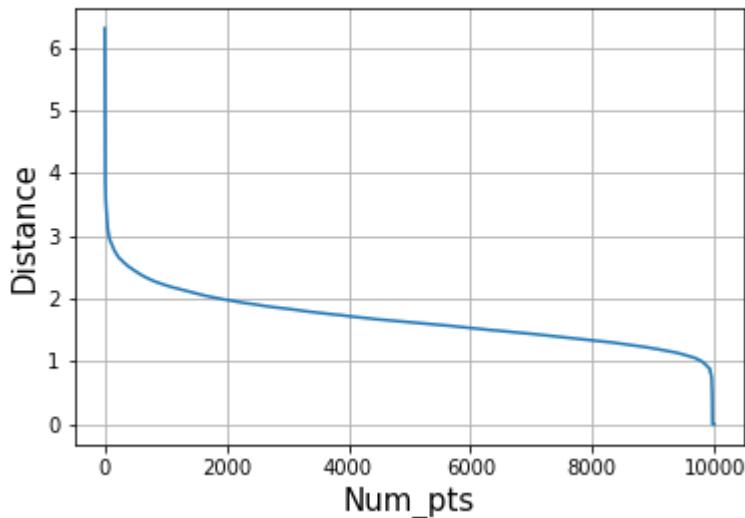


```
1 PrintWordcloud(ppr,tfidf_sent_vect)
```



Plotting Epsilon graph

distance VS Num_pts Plot



eps = 2.4280391942848922

All clusters are [-1 0]

word cloud for cluster -1



word cloud for cluster 0





1. word cloud -1 represents general Noise.
 2. word cloud 2 represents good review about coffee and tea.

▼ [6] Conclusions

```
1 from prettytable import PrettyTable
2 ptable1 = PrettyTable()
3
4 vec1 = "BOW"
5 vec2 = "TFIDF"
6 vec3 = "AVG-W2V"
7 vec4 = "TFIDF-W2V"
8
9 model1 = 'Kmeans Clustering'
10 model2 = 'Agglomerative clustering'
11 model3= 'DBSCAN Clustering'
12
13
14 ptable1.field_names = ["Vectorizer", "Model", "Clusters"]
15 ptable1.add_row([vec1,model1,4])
16 ptable1.add_row([vec2,model1,2])
17 ptable1.add_row([vec3,model1,2])
18 ptable1.add_row([vec4,model1,2])
19 print(ptable1)
20
21 ptable2 = PrettyTable()
22 ptable2.field_names = ["Vectorizer", "Model", "Optimal_Cluster"]
23 ptable2.add_row([vec3,model2,2])
24 ptable2.add_row([vec3,model2,5])
25 ptable2.add_row([vec4,model2,2])
26 ptable2.add_row([vec4,model2,5])
27 print(ptable2)
```

```
28  
29 ptable3 = PrettyTable()  
30 ptable3.field_names = ["Vectorizer", "Model", "Min-pts", "Eps"]  
31 ptable3.add_row([vec3, model3, 50, 2.594])  
32 ptable3.add_row([vec4, model3, 50, 2.428])  
33 print(ptable3)
```



Vectorizer	Model	Clusters	
BOW	Kmeans Clustering	4	
TFIDF	Kmeans Clustering	2	
AVG-W2V	Kmeans Clustering	2	
TFIDF-W2V	Kmeans Clustering	2	

Vectorizer	Model	Optimal_Cluster	
AVG-W2V	Agglomerative clustering	2	
AVG-W2V	Agglomerative clustering	5	
TFIDF-W2V	Agglomerative clustering	2	
TFIDF-W2V	Agglomerative clustering	5	

Vectorizer	Model	Min-pts	Eps	
AVG-W2V	DBSCAN Clustering	50	2.594	
TFIDF-W2V	DBSCAN Clustering	50	2.428	