

# Amazon Fine Food Reviews Analysis

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

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

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

Number of reviews: 568,454

Number of users: 256,059

Number of products: 74,258

Timespan: Oct 1999 - Oct 2012

Number of Attributes/Columns in data: 10

Attribute Information:

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

**Objective:**

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

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

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

## [1]. Reading Data

### [1.1] Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

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

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

```
In [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('C:/Users/Excel/Desktop/vins/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
# 000 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 5000""", 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 (5000, 10)

Out[2]:

	<b>Id</b>	<b>ProductId</b>	<b>UserId</b>	<b>ProfileName</b>	<b>HelpfulnessNumerator</b>	<b>HelpfulnessDenominator</b>
<b>0</b>	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1

	<b>Id</b>	<b>ProductId</b>	<b>UserId</b>	<b>ProfileName</b>	<b>HelpfulnessNumerator</b>	<b>HelpfulnessDenominator</b>
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0
2	3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1

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

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

(80668, 7)

Out[4]:

	<b>UserId</b>	<b>ProductId</b>	<b>ProfileName</b>	<b>Time</b>	<b>Score</b>	<b>Text</b>	<b>COUNT</b>

	UserId	ProductId	ProfileName	Time	Score	Text	COU
0	#oc-R115TNMSPFT9I7	B007Y59HVM	Breyton	1331510400	2	Overall its just OK when considering the price...	2
1	#oc-R11D9D7SHXIJB9	B005HG9ET0	Louis E. Emory "hoppy"	1342396800	5	My wife has recurring extreme muscle spasms, u...	3
2	#oc-R11DNU2NBKQ23Z	B007Y59HVM	Kim Cieszykowski	1348531200	1	This coffee is horrible and unfortunately not ...	2
3	#oc-R11O5J5ZVQE25C	B005HG9ET0	Penguin Chick	1346889600	5	This will be the bottle that you grab from the...	3
4	#oc-R12KPBDL2B5ZD	B007OSBE1U	Christopher P. Presta	1348617600	1	I didnt like this coffee. Instead of telling y...	2



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

Out[5]:

	UserId	ProductId	ProfileName	Time	Score	Text	COU
0	#oc-R115TNMSPFT9I7	B007Y59HVM	Breyton	1331510400	2	Overall its just OK when considering the price...	2

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

	<b>Id</b>	<b>ProductId</b>	<b>UserId</b>	<b>ProfileName</b>	<b>HelpfulnessNumerator</b>	<b>HelpfulnessDenominator</b>
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	2
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	2
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	2
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	2
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan	2	2

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

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

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

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

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

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

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

```
Out[9]: (4986, 10)
```

```
In [11]: #Checking to see how much % of data still remains
(final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

```
Out[11]: 99.72
```

**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 calcualtions

```
In [12]: display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND Id=44737 OR Id=64422
ORDER BY ProductID
""", con)

display.head()
```

Out[12]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	Helpfuln
0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	1
1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	2

```
In [13]: final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [14]: #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()
```

```
(4986, 10)
```

```
Out[14]: 1    4178  
0     808  
Name: Score, dtype: int64
```

```
In [15]: final['Time']=pd.to_datetime(final['Time'],unit='s')  
final=final.sort_values(by='Time')
```

## [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 [16]: # 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)
```

This was a really good idea and the final product is outstanding. I use the decals on my car window and everybody asks where i bought the decal s i made. Two thumbs up!

=====

These are thin,crisp, fragrant cookies and are very delicious and tast y. They are excellent with a glass of cold almond milk or hot herbal te a. (my choices) If you like ginger snaps you will love Lars ginger snap s.

=====

Green Mountain "Nantucket Blend" K-Cups make a very good cup of coffee in my <a href="http://www.amazon.com/gp/product/B000AQPMHA">Keurig B-40 B40 Elite Gourmet Single-Cup Home-Brewing System</a>. This is a very sm ooth tasting brew that my wife prefers over the <a href="http://www.ama zon.com/gp/product/B0029XDZIK">Coffee People, Donut Shop K-Cups for Keu rig Brewers (Pack of 50) [Amazon Frustration-Free Packaging</a>] I gene rally drink in the morning.<br /><br />These are good on both "Small" a nd "Large" cup settings as well.<br /><br />Highly Recommended!<br /><b r />CFH

=====

Besides being smaller than runts, they look the same and have the same consistency. Unfortunately, they taste nothing like banana runts...nor do they even taste good. Yucky stuff. Trying to return with vendor.

```
In [17]: # remove urls from text python: https://stackoverflow.com/a/40823105/40  
84039  
sent_0 = re.sub(r"http\S+", "", sent_0)  
sent_1000 = re.sub(r"http\S+", "", sent_1000)  
sent_150 = re.sub(r"http\S+", "", sent_1500)  
sent_4900 = re.sub(r"http\S+", "", sent_4900)  
  
print(sent_0)
```

This was a really good idea and the final product is outstanding. I use the decals on my car window and everybody asks where i bought the decal s i made. Two thumbs up!

```
In [18]: # 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)
```

This was a really good idea and the final product is outstanding. I use the decals on my car window and everybody asks where i bought the decal s i made. Two thumbs up!

=====

These are thin, crisp, fragrant cookies and are very delicious and tasty. They are excellent with a glass of cold almond milk or hot herbal tea. (my choices) If you like ginger snaps you will love Lars ginger snaps.

=====

Green Mountain "Nantucket Blend" K-Cups make a very good cup of coffee in my Keurig B-40 B40 Elite Gourmet Single-Cup Home-Brewing System. This is a very smooth tasting brew that my wife prefers over the Coffee People, Donut Shop K-Cups for Keurig Brewers (Pack of 50) [Amazon Frustration-Free Packaging] I generally drink in the morning. These are good on both "Small" and "Large" cup settings as well. Highly Recommended! CFH

=====

Besides being smaller than runts, they look the same and have the same consistency. Unfortunately, they taste nothing like banana runts...nor do they even taste good. Yucky stuff. Trying to return with vendor.

```
In [19]: # 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 [20]: sent_1500 = decontracted(sent_1500)
print(sent_1500)
print("*50)
```

Green Mountain "Nantucket Blend" K-Cups make a very good cup of coffee in my <a href="http://www.amazon.com/gp/product/B000AQPMHA">Keurig B-40 B40 Elite Gourmet Single-Cup Home-Brewing System</a>. This is a very smooth tasting brew that my wife prefers over the <a href="http://www.amazon.com/gp/product/B0029XDZIK">Coffee People, Donut Shop K-Cups for Keurig Brewers (Pack of 50) [Amazon Frustration-Free Packaging</a>] I generally drink in the morning.<br /><br />These are good on both "Small" and "Large" cup settings as well.<br /><br />Highly Recommended!<br /><b>r />CFH

=====

```
In [21]: #remove words with numbers python: https://stackoverflow.com/a/18082370/4084039
sent_0 = re.sub("\S*\d\S*", "", sent_0).strip()
print(sent_0)
```

This was a really good idea and the final product is outstanding. I use the decals on my car window and everybody asks where i bought the decal s i made. Two thumbs up!

```
In [22]: #remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)
print(sent_1500)
```

Green Mountain Nantucket Blend K Cups make a very good cup of coffee in my a href http www amazon com gp product B000AQPMHA Keurig B 40 B40 Elite Gourmet Single Cup Home Brewing System a This is a very smooth tasting brew that my wife prefers over the a href http www amazon com gp product B0029XDZIK Coffee People Donut Shop K Cups for Keurig Brewers Pack of 50 Amazon Frustration Free Packaging a I generally drink in the morn ing br br These are good on both Small and Large cup settings as well b r br Highly Recommended br br CFH

```
In [23]: # https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
# <br /><br /> ==> after the above steps, we are getting "br br"
# we are including them into stop words list
# instead of <br /> if we have <br/> these tags would have removed in
```

*the 1st step*

```
stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'o
urs', '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", 'it
s', 'itself', 'they', 'them', 'their', \
    'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'th
is', 'that', "that'll", 'these', 'those', \
    'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'h
ave', 'has', 'had', 'having', 'do', 'does', \
    'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or',
'because', 'as', 'until', 'while', 'of', \
    'at', 'by', 'for', 'with', 'about', 'against', 'between',
'into', 'through', 'during', 'before', 'after', \
    'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out',
'on', 'off', 'over', 'under', 'again', 'further', \
    'then', 'once', 'here', 'there', 'when', 'where', 'why', 'h
ow', 'all', 'any', 'both', 'each', 'few', 'more', \
    'most', 'other', 'some', 'such', 'only', 'own', 'same', 's
o', 'than', 'too', 'very', \
    's', 't', 'can', 'will', 'just', 'don', "don't", 'should',
"should've", 'now', 'd', 'll', 'm', 'o', 're', \
    've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't",
'didn', "didn't", 'doesn', "doesn't", 'hadn', \
    "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "is
n't", 'ma', 'mightn', "mightn't", 'mustn', \
    "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn',
"shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
    'won', "won't", 'wouldn', "wouldn't"])
```

In [24]: # 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()
```

```
sentance = decontracted(sentance)
sentance = re.sub("\S*\d\S*", "", sentance).strip()
sentance = re.sub('[^A-Za-z]+', ' ', sentance)
# https://gist.github.com/sebleier/554280
sentance = ' '.join(e.lower() for e in sentance.split() if e.lower()
() not in stopwords)
preprocessed_reviews.append(sentance.strip())
```

```
100%|██████████| 4986/4986 [00:02<00:00, 200
9.56it/s]
```

In [25]: preprocessed\_reviews[1500]

Out[25]: 'green mountain nantucket blend k cups make good cup coffee generally d  
rink morning good small large cup settings well highly recommended cfh'

## [3.2] Preprocessing Review Summary

In [26]: final["preprocessed\_reviews"] = preprocessed\_reviews

In [28]: X=final['preprocessed\_reviews'].values  
X.shape

Out[28]: (4986,)

```
In [29]: sent_x = []
for sent in X :
    sent_x.append(sent.split())

# Train your own Word2Vec model using your own train text corpus
# min_count = 5 considers only words that occurred atleast 5 times
w2v_model=Word2Vec(sent_x,min_count=5,size=50, workers=4)

w2v_words = list(w2v_model.wv.vocab)
print("number of words that occurred minimum 5 times ",len(w2v_words))
```

```

# compute average word2vec for each review for sent_x .
train_vectors = []
for sent in tqdm(sent_x):
    sent_vec = np.zeros(50)
    cnt_words = 0
    for word in sent: #
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    train_vectors.append(sent_vec)

data = train_vectors

```

number of words that occurred minimum 5 times 3817

100%|██████████| 4986/4986 [00:05<00:00, 91  
2.30it/s]

In [33]: `from sklearn.cluster import AgglomerativeClustering  
agglo=AgglomerativeClustering(n_clusters=2).fit(data)`

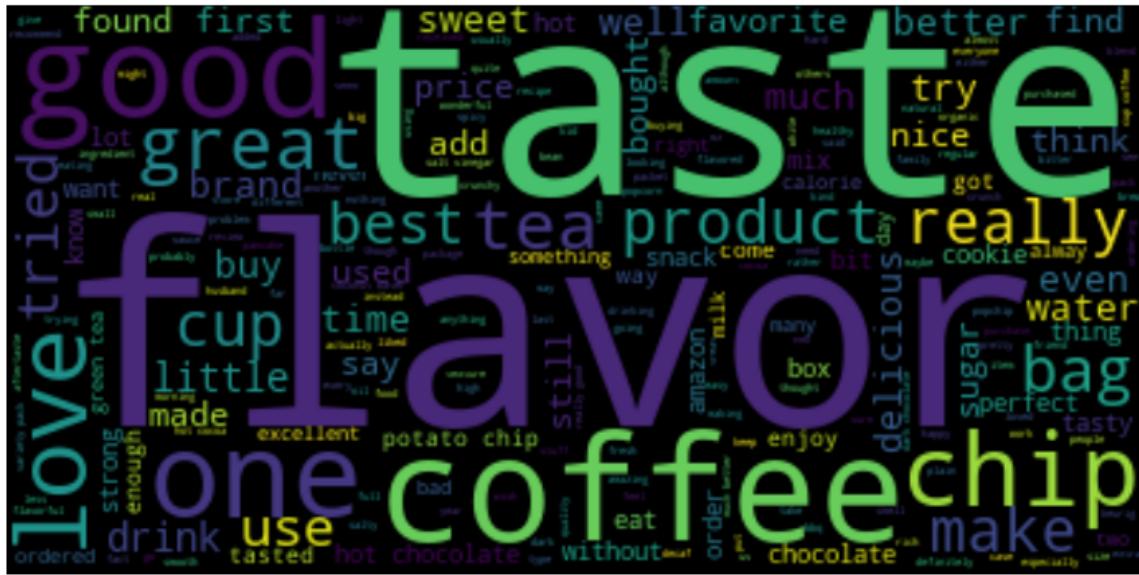
In [34]: `cluster1 = []  
cluster2 = []  
  
for i in range(agglo.labels_.shape[0]):  
 if agglo.labels_[i] == 0:  
 cluster1.append(preprocessed_reviews[i])  
 else :  
 cluster2.append(preprocessed_reviews[i])  
print("No. of reviews in Cluster-1 : ",len(cluster1))  
print("\nNo. of reviews in Cluster-2 : ",len(cluster2))`

No. of reviews in Cluster-1 : 2603

No. of reviews in Cluster-2 : 2383

```
In [37]: text_1=cluster1  
text_2=cluster2  
  
lst=[text_1,text_2]  
for i in lst:  
    from wordcloud import WordCloud, STOPWORDS  
    stopwords = set(STOPWORDS)  
  
    wordcloud = WordCloud(max_words=1000).generate(str(i))  
    plt.figure(figsize = (10, 10), facecolor = None)  
    plt.imshow(wordcloud, interpolation="bilinear")  
    plt.axis("off")  
    plt.tight_layout(pad = 0)  
    plt.show()
```





1. cluster1: good,great,product,food,one,taste.
  2. cluster2: taste ,flavour,coffee,good.

```
In [38]: agglo=AgglomerativeClustering(n_clusters=3).fit(data)
```

```
In [40]: cluster1 = []
          cluster2 = []
          cluster3 = []
```

```
for i in range(agglo.labels_.shape[0]):
    if agglo.labels_[i] == 0:
        cluster1.append(preprocessed_reviews[i])
    elif agglo.labels_[i] == 1:
        cluster2.append(preprocessed_reviews[i])
    else :
        cluster3.append(preprocessed_reviews[i])
print("No. of reviews in Cluster-1 : ",len(cluster1))
```

```
print("\nNo. of reviews in Cluster-1 : ",len(cluster1))
print("\nNo. of reviews in Cluster-2 : ",len(cluster2))
print("\nNo. of reviews in Cluster-3 : ",len(cluster3))

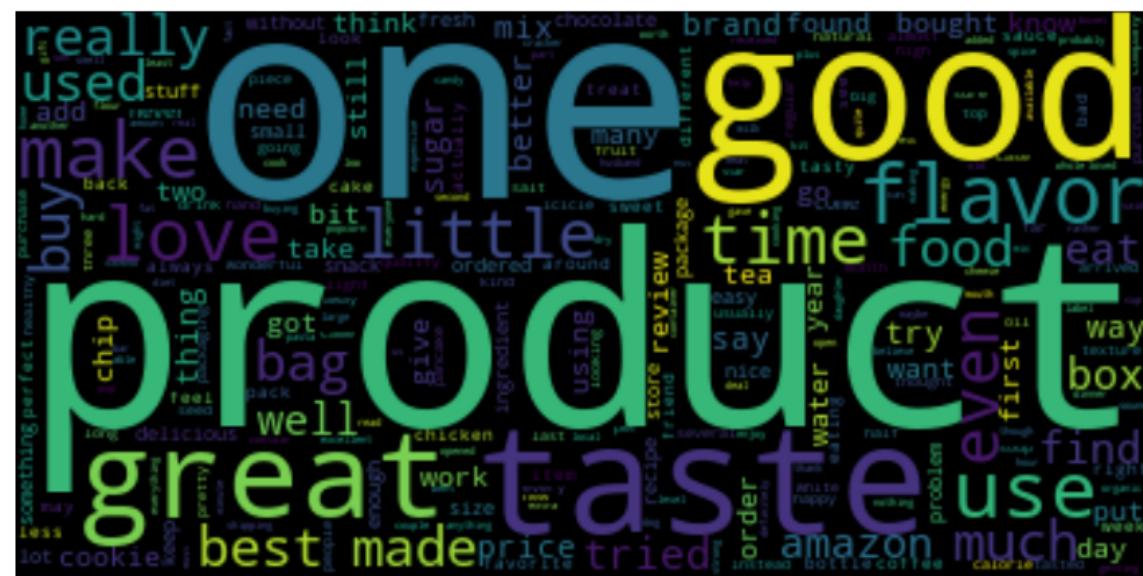
No. of reviews in Cluster-1 : 2383

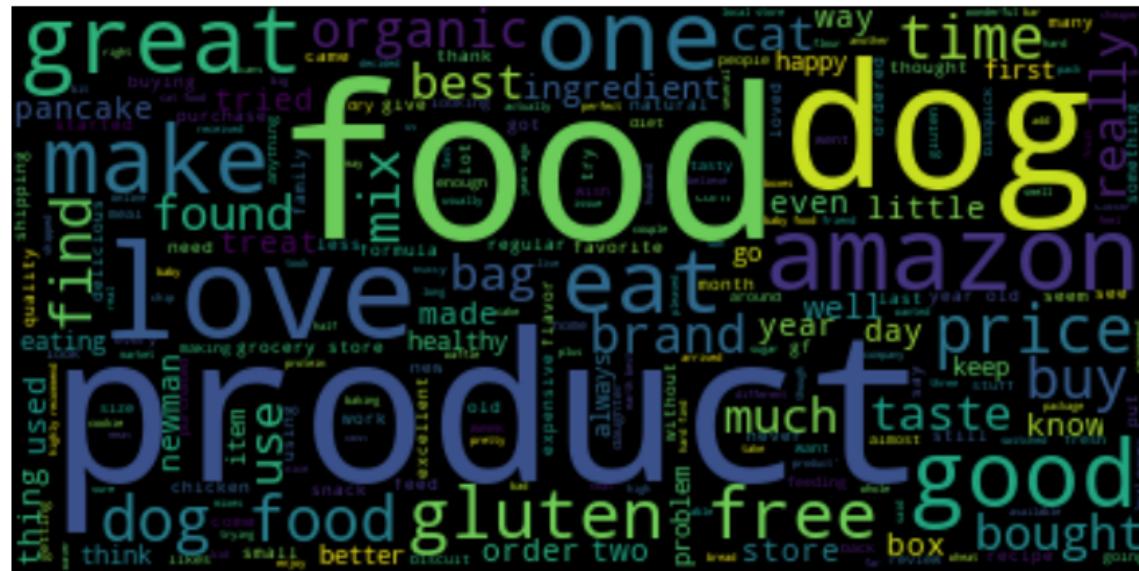
No. of reviews in Cluster-2 : 1362

No. of reviews in Cluster-3 : 1241
```

```
In [41]: text_1=cluster1
text_2=cluster2
text_3=cluster3
lst=[text_1,text_2,text_3]
for i in lst:
    from wordcloud import WordCloud, STOPWORDS
    stopwords = set(STOPWORDS)

    wordcloud = WordCloud(max_words=1000).generate(str(i))
    plt.figure(figsize = (10, 10), facecolor = None)
    plt.imshow(wordcloud, interpolation="bilinear")
    plt.axis("off")
    plt.tight_layout(pad = 0)
    plt.show()
```





1. cluster1: good,coffee,taste, flavor,one,tea.
  2. cluster2: one,good,product,great,taste.
  3. cluster3: great,food,dog,product,good.

```
In [43]: aggro=AgglomerativeClustering(n_clusters=4).fit(data)
```

```
In [44]: cluster1 = []
        cluster2 = []
        cluster3 = []
        cluster4 = []

        for i in range(agglo.labels_.shape[0]):
            if agglo.labels_[i] == 0:
                cluster1.append(preprocessed_reviews[i])
            elif agglo.labels_[i] == 1:
```

```
        cluster2.append(preprocessed_reviews[i])
    elif agglo.labels_[i] == 2:
        cluster3.append(preprocessed_reviews[i])
    else :
        cluster4.append(preprocessed_reviews[i])
print("No. of reviews in Cluster-1 : ",len(cluster1))
print("\nNo. of reviews in Cluster-2 : ",len(cluster2))
print("\nNo. of reviews in Cluster-3 : ",len(cluster3))
print("\nNo. of reviews in Cluster-4 : ",len(cluster4))
```

```
No. of reviews in Cluster-1 : 1362
```

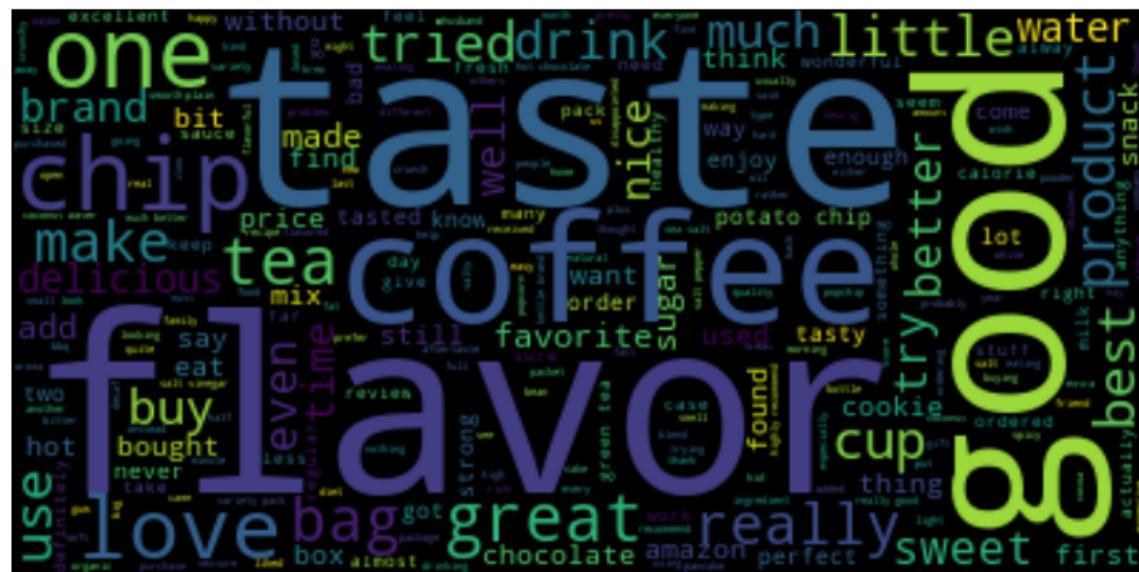
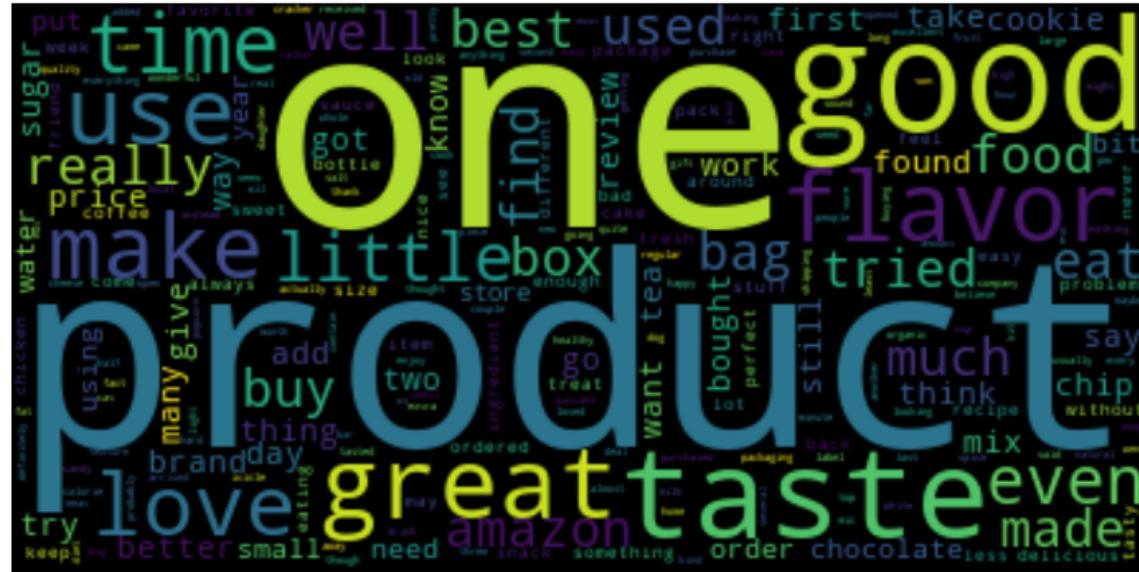
```
No. of reviews in Cluster-2 : 2162
```

```
No. of reviews in Cluster-3 : 1241
```

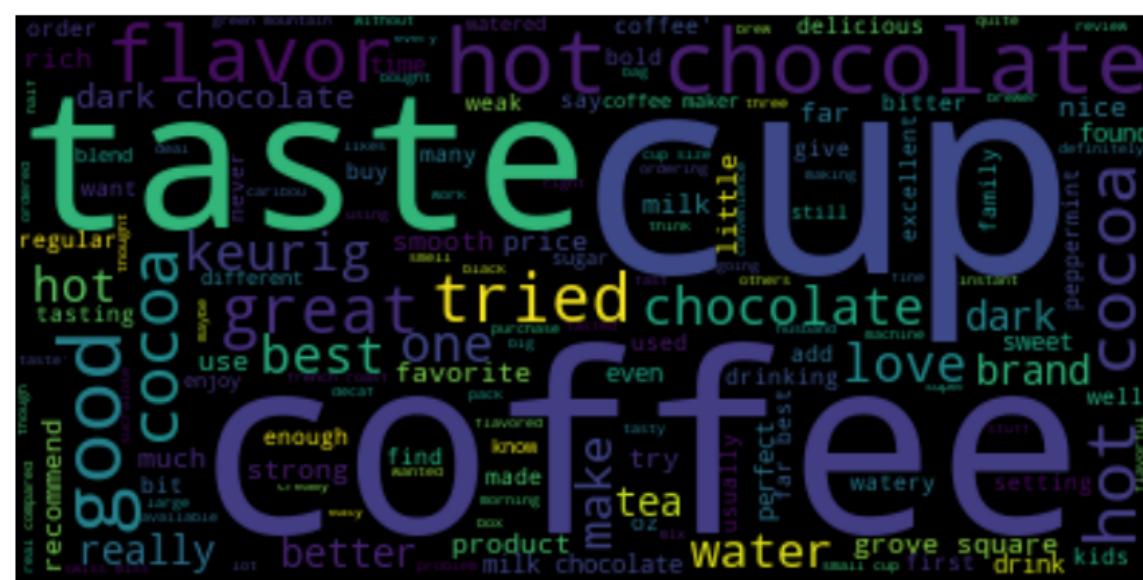
```
No. of reviews in Cluster-4 : 221
```

```
In [46]: text_1=cluster1
text_2=cluster2
text_3=cluster3
text_4=cluster4
lst=[text_1,text_2,text_3,text_4]
for i in lst:
    from wordcloud import WordCloud, STOPWORDS
    stopwords = set(STOPWORDS)
```

```
wordcloud = WordCloud(max_words=1000).generate(str(i))
plt.figure(figsize = (10, 10), facecolor = None)
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis("off")
plt.tight_layout(pad = 0)
plt.show()
```



gluten\_free newman\_starters little\_year\_old\_pets one\_time\_a\_day



1. cluster1: time, one,good,product,great,taste.
  2. cluster2: one,taste,coffee,flavor,good,great.

3. cluster3: dog,food,amazon,good,love.

4. cluster4: cup,coffee,tried, chocolate.

```
In [56]: agglo=AgglomerativeClustering(n_clusters=5).fit(data)
```

```
In [58]: cluster1 = []
cluster2 = []
cluster3 = []
cluster4 = []
cluster5 = []
```

```
for i in range(agglo.labels_.shape[0]):
    if agglo.labels_[i] == 0:
        cluster1.append(preprocessed_reviews[i])
    elif agglo.labels_[i] == 1:
        cluster2.append(preprocessed_reviews[i])
    elif agglo.labels_[i] == 2:
        cluster3.append(preprocessed_reviews[i])
    elif agglo.labels_[i] == 3:
        cluster4.append(preprocessed_reviews[i])
    else :
        cluster5.append(preprocessed_reviews[i])
print("No. of reviews in Cluster-1 : ",len(cluster1))
print("\nNo. of reviews in Cluster-2 : ",len(cluster2))
print("\nNo. of reviews in Cluster-3 : ",len(cluster3))
print("\nNo. of reviews in Cluster-4 : ",len(cluster4))
print("\nNo. of reviews in Cluster-5 : ",len(cluster5))
```

No. of reviews in Cluster-1 : 2162

No. of reviews in Cluster-2 : 1347

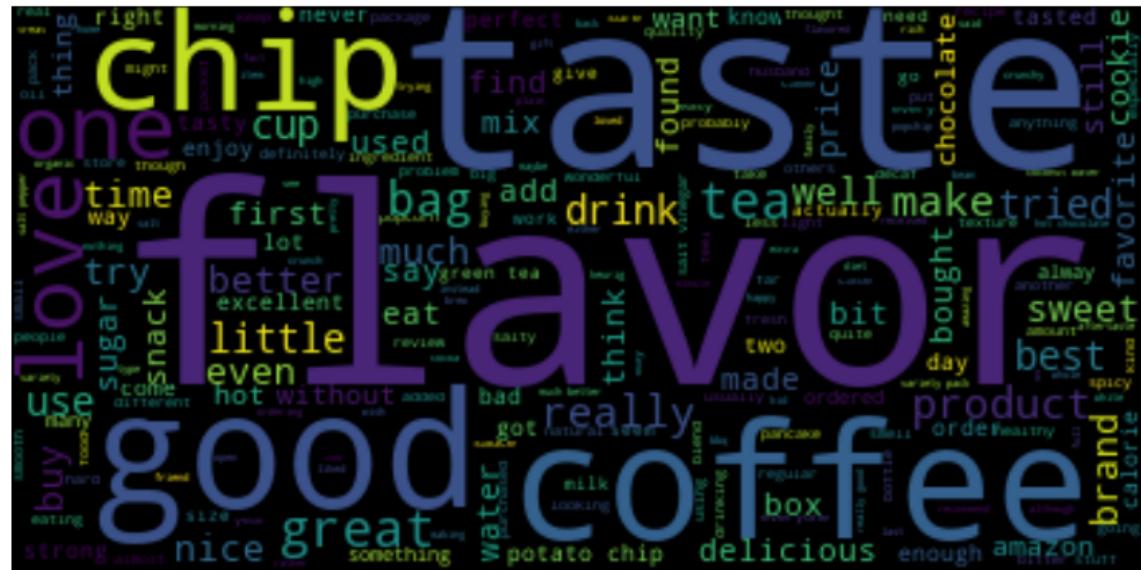
No. of reviews in Cluster-3 : 1241

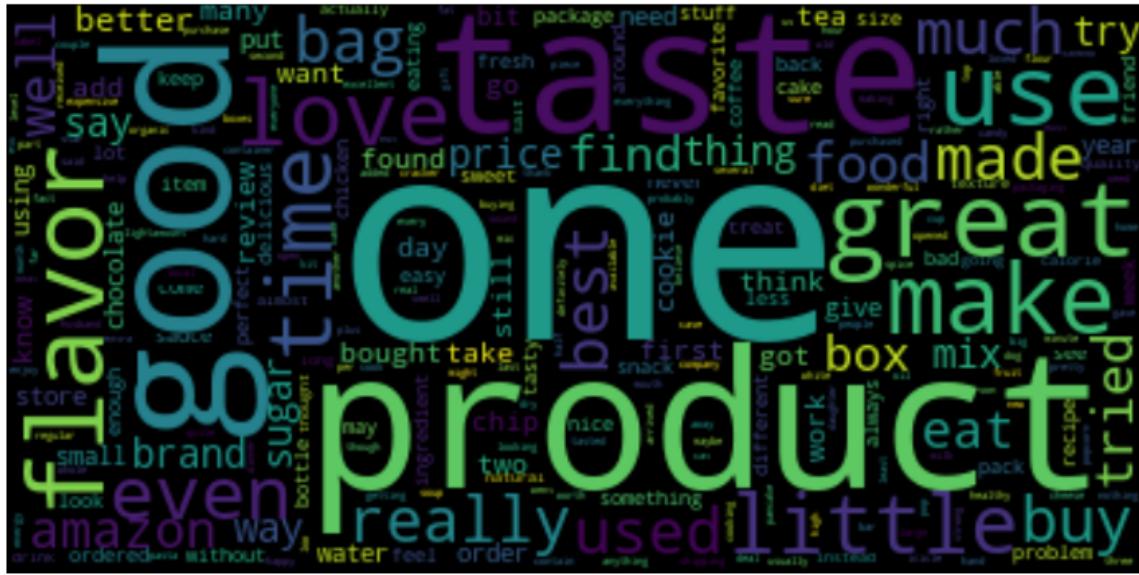
No. of reviews in Cluster-4 : 221

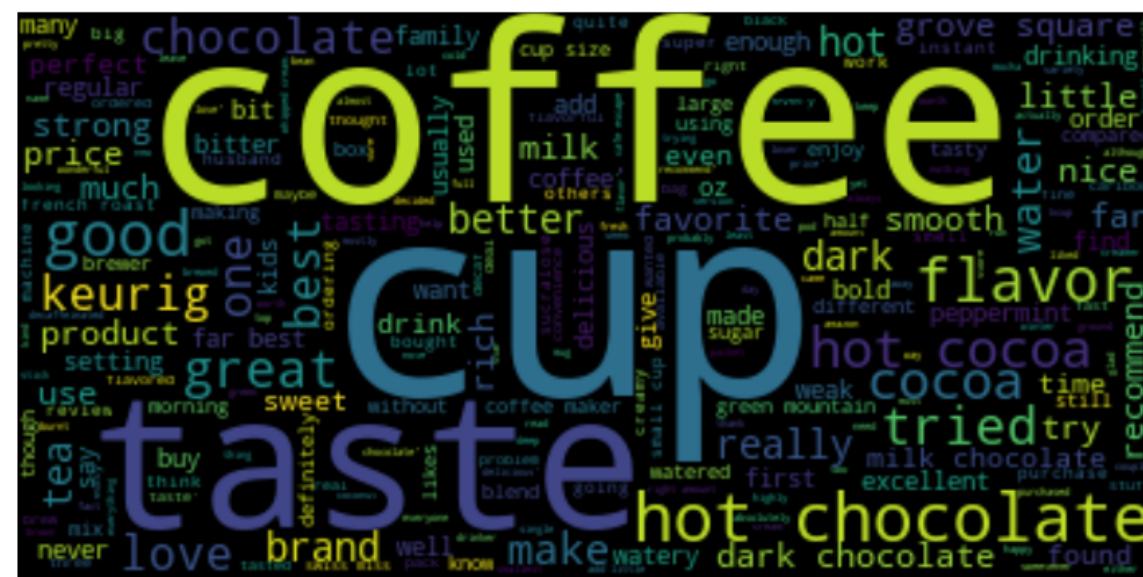
No. of reviews in Cluster-5 : 15

```
In [59]: text_1=cluster1  
text_2=cluster2  
text_3=cluster3  
text_4=cluster4  
lst=[text_1,text_2,text_3,text_4]  
for i in lst:  
    from wordcloud import WordCloud, STOPWORDS  
    stopwords = set(STOPWORDS)
```

```
wordcloud = WordCloud(max_words=1000).generate(str(i))
plt.figure(figsize = (10, 10), facecolor = None)
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis("off")
plt.tight_layout(pad = 0)
plt.show()
```







1. cluster1: chip,taste,flavor,good,coffee.
  2. cluster2: flavor,good,one product,great,make.
  3. cluster3: product,amazon,food,dog,eat,make.
  4. cluster4: coffee,cup,taste,hot chocolate,brand.
  5. cluster5: chip,product,time,one good.

## TF-IDF W2V

```
In [65]: model = TfidfVectorizer()
model.fit(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_)))
# 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 (sent_x): # 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 values 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
data=tfidf_sent_vectors
```

```
In [66]: agglo_two=AgglomerativeClustering(n_clusters=2).fit(data)
```

```
In [68]: cluster1 = []
cluster2 = []

for i in range(agglo.labels_.shape[0]):
    if agglo.labels_[i] == 0:
        cluster1.append(preprocessed_reviews[i])
    else :
        cluster2.append(preprocessed_reviews[i])
print("No. of reviews in Cluster-1 : ",len(cluster1))
print("\nNo. of reviews in Cluster-2 : ",len(cluster2))
```

No. of reviews in Cluster-1 : 2162

No. of reviews in Cluster-2 : 2824

```
In [69]: text_1=cluster1
text_2=cluster2

lst=[text_1,text_2]
for i in lst:
    from wordcloud import WordCloud, STOPWORDS
    stopwords = set(STOPWORDS)

    wordcloud = WordCloud(max_words=1000).generate(str(i))
    plt.figure(figsize = (10, 10), facecolor = None)
    plt.imshow(wordcloud, interpolation="bilinear")
    plt.axis("off")
    plt.tight_layout(pad = 0)
    plt.show()
```



1. cluster1: coffee,good,flavor,taste,chip,tea.
2. cluster2: product,one,food,great,good.

```
In [70]: agglo=AgglomerativeClustering(n_clusters=3).fit(data)
```

```
In [72]: cluster1 = []
cluster2 = []
cluster3 = []

for i in range(agglo.labels_.shape[0]):
    if agglo.labels_[i] == 0:
        cluster1.append(preprocessed_reviews[i])
    elif agglo.labels_[i] == 1:
        cluster2.append(preprocessed_reviews[i])
    else :
        cluster3.append(preprocessed_reviews[i])
print("No. of reviews in Cluster-1 : ",len(cluster1))
print("\nNo. of reviews in Cluster-2 : ",len(cluster2))
print("\nNo. of reviews in Cluster-3 : ",len(cluster3))
```

```
No. of reviews in Cluster-1 : 2308
```

```
No. of reviews in Cluster-2 : 1152
```

```
No. of reviews in Cluster-3 : 1526
```

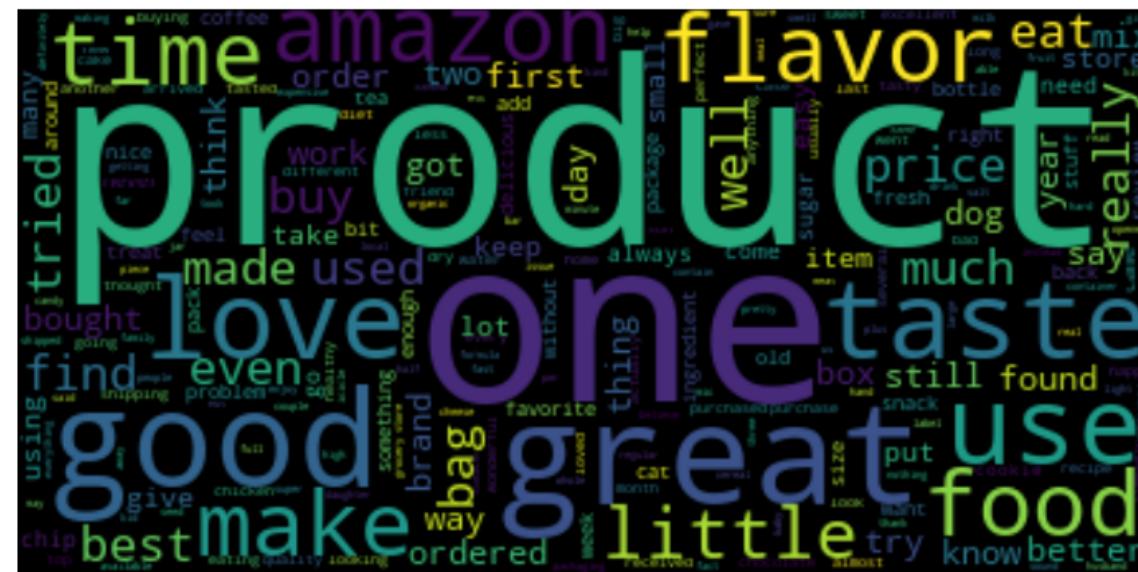
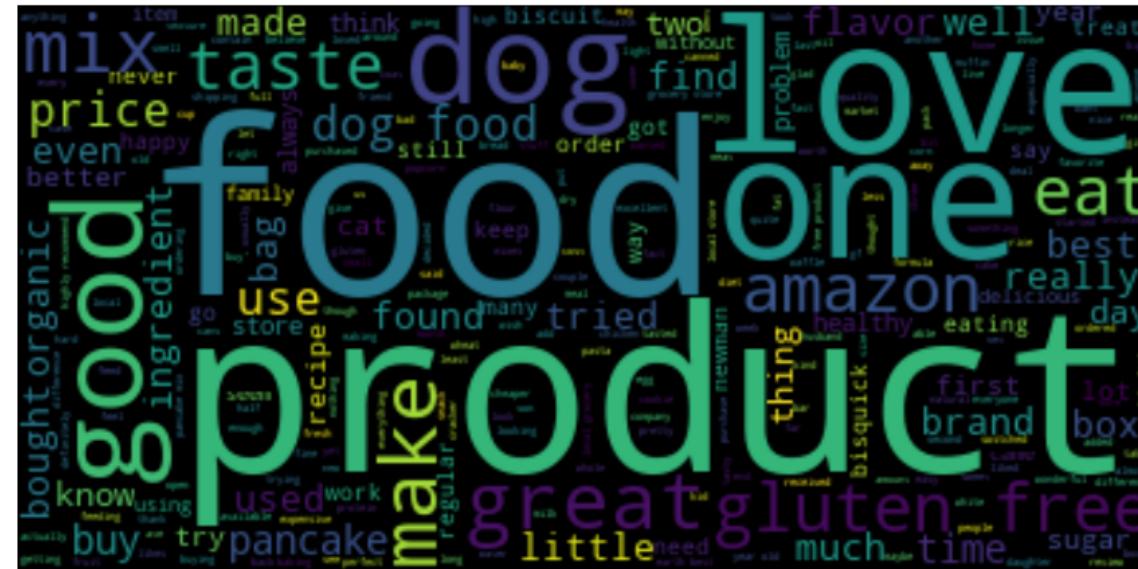
```
In [74]: text_1=cluster1
text_2=cluster2
text_3=cluster3

lst=[text_1,text_2,text_3]
for i in lst:
    from wordcloud import WordCloud, STOPWORDS
    stopwords = set(STOPWORDS)

    wordcloud = WordCloud(max_words=1000).generate(str(i))
    plt.figure(figsize = (10, 10), facecolor = None)
```

```
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis("off")
plt.tight_layout(pad = 0)
plt.show()
```





1. cluster1:flavor,taste,bag,coffee,love,tea,drink.
2. cluster2:taste,dog,love,food,one,product.
3. cluster3:product,love,one,taste,great,usefull,little.

```
In [81]: agglo=AgglomerativeClustering(n_clusters=4).fit(data)
```

```
In [82]: cluster1 = []
cluster2 = []
cluster3 = []
cluster4 = []

for i in range(agglo.labels_.shape[0]):
    if agglo.labels_[i] == 0:
        cluster1.append(preprocessed_reviews[i])
    elif agglo.labels_[i] == 1:
        cluster2.append(preprocessed_reviews[i])
    elif agglo.labels_[i] == 2:
        cluster3.append(preprocessed_reviews[i])
    else :
        cluster4.append(preprocessed_reviews[i])
print("No. of reviews in Cluster-1 : ",len(cluster1))
print("\nNo. of reviews in Cluster-2 : ",len(cluster2))
print("\nNo. of reviews in Cluster-3 : ",len(cluster3))
print("\nNo. of reviews in Cluster-4 : ",len(cluster4))
```

No. of reviews in Cluster-1 : 1152

No. of reviews in Cluster-2 : 1830

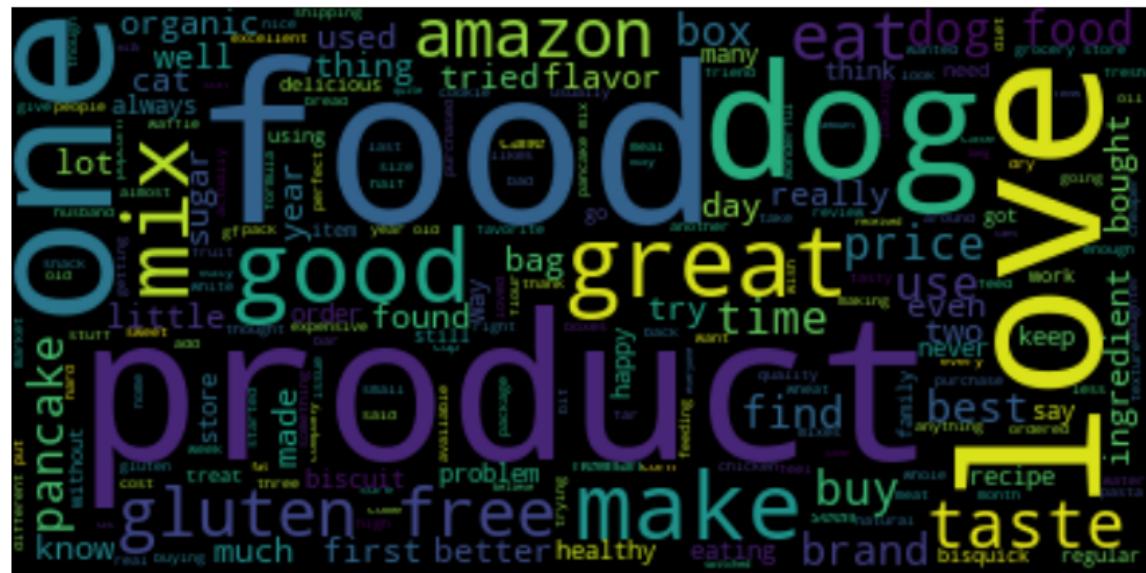
No. of reviews in Cluster-3 : 1526

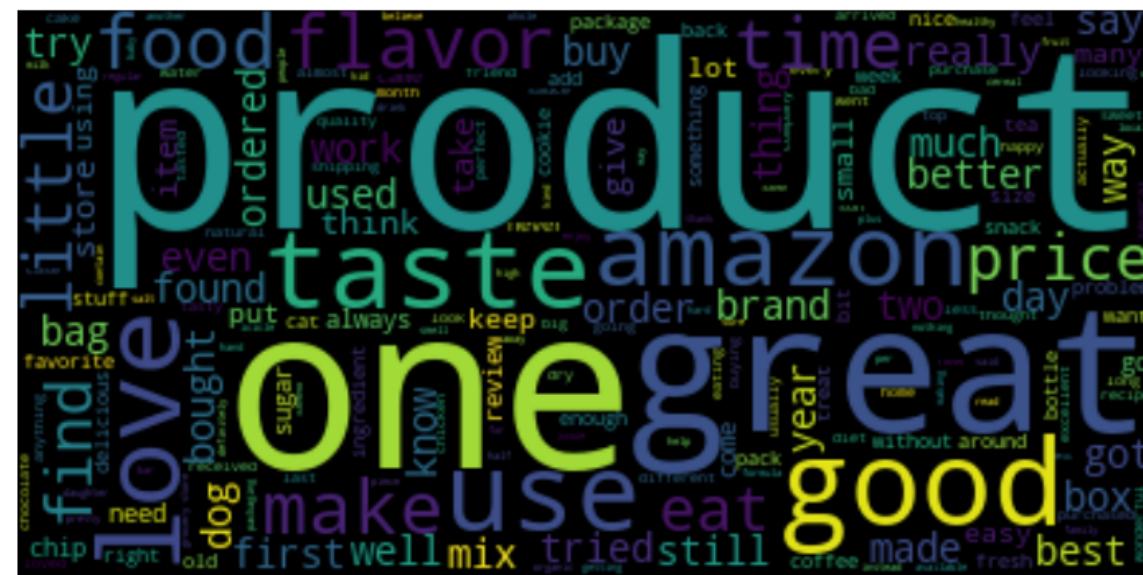
No. of reviews in Cluster-4 : 478

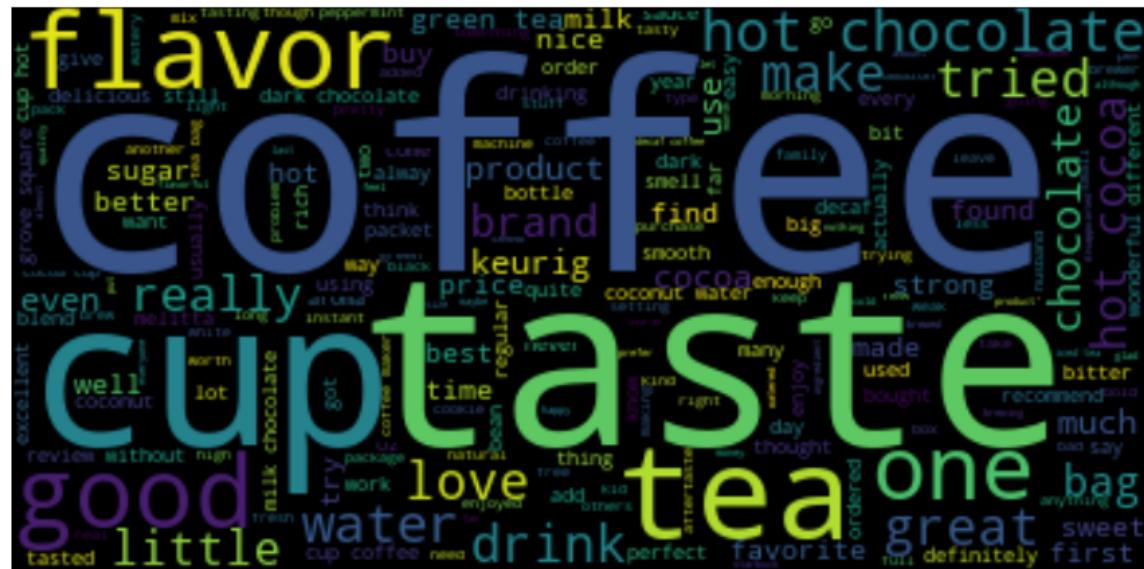
```
In [83]: text_1=cluster1
text_2=cluster2
text_3=cluster3
text_4=cluster4
```

```
lst=[text_1,text_2,text_3,text_4]
for i in lst:
    from wordcloud import WordCloud, STOPWORDS
    stopwords = set(STOPWORDS)

    wordcloud = WordCloud(max_words=1000).generate(str(i))
    plt.figure(figsize = (10, 10), facecolor = None)
    plt.imshow(wordcloud, interpolation="bilinear")
    plt.axis("off")
    plt.tight_layout(pad = 0)
    plt.show()
```







1. cluster1:food,dog,love,product,make,buy.
  2. cluster2:flavor,taste,good,chip,coffee,one,love.
  3. cluster3:product,taste,amazon,great,good,love.
  4. cluster4:flavor.coffee.cup.taste.tea.one.

```
In [84]: agglo=AgglomerativeClustering(n_clusters=5).fit(data)
```

```
In [85]: cluster1 = []
          cluster2 = []
```

```
cluster3 = []
cluster4 = []
cluster5 = []

for i in range(agglo.labels_.shape[0]):
    if agglo.labels_[i] == 0:
        cluster1.append(preprocessed_reviews[i])
    elif agglo.labels_[i] == 1:
        cluster2.append(preprocessed_reviews[i])
    elif agglo.labels_[i] == 2:
        cluster3.append(preprocessed_reviews[i])
    elif agglo.labels_[i] == 3:
        cluster4.append(preprocessed_reviews[i])

    else :
        cluster5.append(preprocessed_reviews[i])
print("No. of reviews in Cluster-1 : ",len(cluster1))
print("\nNo. of reviews in Cluster-2 : ",len(cluster2))
print("\nNo. of reviews in Cluster-3 : ",len(cluster3))
print("\nNo. of reviews in Cluster-4 : ",len(cluster4))
print("\nNo. of reviews in Cluster-5 : ",len(cluster5))
```

```
No. of reviews in Cluster-1 : 1830
```

```
No. of reviews in Cluster-2 : 1062
```

```
No. of reviews in Cluster-3 : 1526
```

```
No. of reviews in Cluster-4 : 478
```

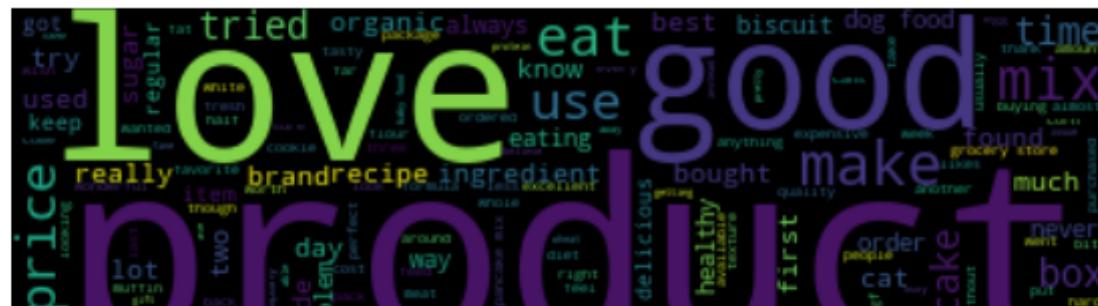
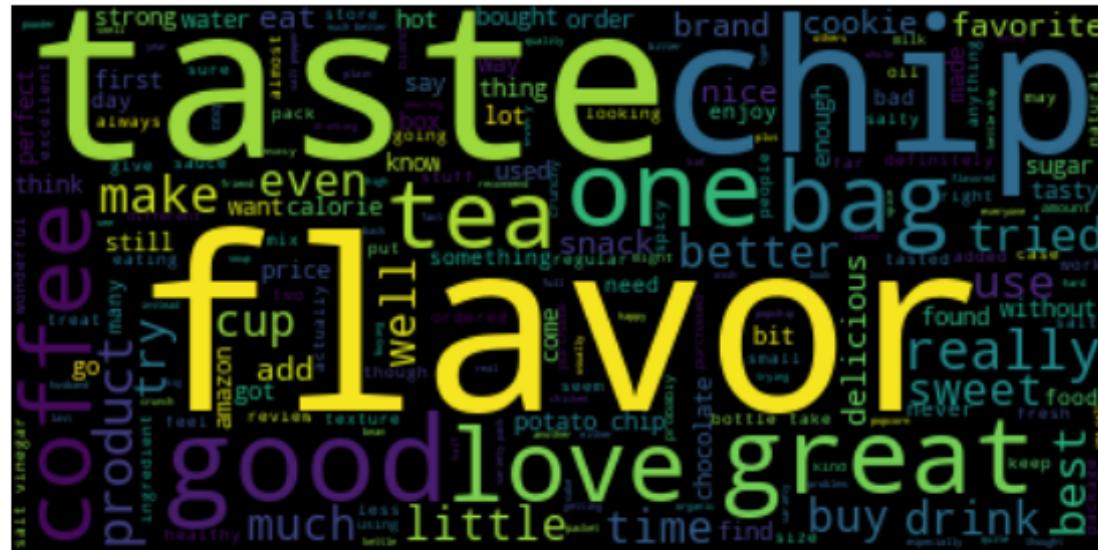
```
No. of reviews in Cluster-5 : 90
```

In [88]:

```
text_1=cluster1
text_2=cluster2
text_3=cluster3
text_4=cluster4
text_5=cluster5

lst=[text_1,text_2,text_3,text_4,text_5]
```

```
for i in lst:  
    from wordcloud import WordCloud, STOPWORDS  
    stopwords = set(STOPWORDS)  
  
    wordcloud = WordCloud(max_words=1000).generate(str(i))  
    plt.figure(figsize = (8, 18), facecolor = None)  
    plt.imshow(wordcloud, interpolation="bilinear")  
    plt.axis("off")  
    plt.tight_layout(pad = 0)  
    plt.show()
```









1. cluster1:taste,chip,love,great,good,coffee.
  2. cluster2:love,good,product,one,food.
  3. cluster3:love,product,one,great,good,amazon,taste.
  4. cluster4:tea,coffee,cup,taste,hot,chocolate.
  5. cluster5:food,dog,gluten,love,good,brand.

## **CONCLUSION**

```
In [3]: from tabulate import tabulate  
print(tabulate ([[ 'AVG-W2V' ,(2,3,4,5)] , [ 'TFIDF-W2V' ,(2,3,4,5)]],  
headers=[ 'Vectorizer' , 'CLUSTER']))
```

Vectorizer	CLUSTER
AVG-W2V	(2, 3, 4, 5)
TFIDF-W2V	(2, 3, 4, 5)

### **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

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

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

### [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

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

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

### [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

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

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

### [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

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

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

### [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

## [5.2] Agglomerative Clustering

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

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

### [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

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

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

### [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
```

## [5.3] DBSCAN Clustering

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

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

### [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
```

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

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

### [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
```

## [6] Conclusions

```
In [4]: # Please compare all your models using Prettytable library.  
# You can have 3 tables, one each for kmeans, agglomerative and dbscan
```

