

# 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

```

In [2]:

```

# using SQLite Table to read data.
con = sqlite3.connect('database.sqlite')

# filtering only positive and negative reviews i.e.
# not taking into consideration those reviews with Score=3
# SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000, will give top 500000 data points
# you can change the number to any other number based on your computing power

# filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000""", co
n)
# for tsne assignment you can take 5k data points

filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3 """, 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 (525814, 10)

Out[2]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	1	1303862400	Good Quality Dog Food
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	0	1346976000	Not as Advertised
2	3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia	1	1	1	1219017600	"Delight" says it all

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary
--	----	-----------	--------	-------------	----------------------	------------------------	-------	------	---------

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
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
```

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

Out[7]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summ
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA VANII WAFE
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA VANII WAFE
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA VANII WAFE
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA VANII WAFE
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACK QUADRA VANII WAFE

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]:

```
(364173, 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]:

69.25890143662969

**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	Score	Time	Summary
0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	1	5	1224892800	Bought This for My Son at College
1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	2	4	1212883200	Pure cocoa taste with crunchy almonds inside

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()
```

(364171, 10)

Out[13]:

```
1    307061
0     57110
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)
```

this witty little book makes my son laugh at loud. i recite it in the car as we're driving along and he always can sing the refrain. he's learned about whales, India, drooping roses: i love all the new words this book introduces and the silliness of it all. this is a classic book i am willing to bet my son will STILL be able to recite from memory when he is in college

=====

I was really looking forward to these pods based on the reviews. Starbucks is good, but I prefer bolder taste.... imagine my surprise when I ordered 2 boxes - both were expired! One expired back in 2005 for gosh sakes. I admit that Amazon agreed to credit me for cost plus part of shipping, but geez, 2 years expired!!! I'm hoping to find local San Diego area shoppe that carries pods so that I can try something different than starbucks.

=====

Great ingredients although, chicken should have been 1st rather than chicken broth, the only thing I do not think belongs in it is Canola oil. Canola or rapeseed is not something a dog would ever find in nature and if it did find rapeseed in nature and eat it, it would poison them. Today's Food industries have convinced the masses that Canola oil is a safe and even better oil than olive or virgin coconut, facts though say otherwise. Until the late 70's it was poisonous until they figured out a way to fix that. I still like it but it could be better.

=====

Can't do sugar. Have tried scores of SF Syrups. NONE of them can touch the excellence of this product.<br /><br />Thick, delicious. Perfect. 3 ingredients: Water, Maltitol, Natural Maple Flavor. PERIOD. No chemicals. No garbage.<br /><br />Have numerous friends & family members hooked on this stuff. My husband & son, who do NOT like "sugar free" prefer this over major label regular syrup.<br /><br />I use this as my SWEETENER in baking: cheesecakes, white brownies, muffins, pumpkin pies, etc... Unbelievably delicious...<br /><br />Can you tell I like it? :)

=====

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)
```

this witty little book makes my son laugh at loud. i recite it in the car as we're driving along and he always can sing the refrain. he's learned about whales, India, drooping roses: i love all the new words this book introduces and the silliness of it all. this is a classic book i am willing to bet my son will STILL be able to recite from memory when he is in college

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)

```

this witty little book makes my son laugh at loud. i recite it in the car as we're driving along and he always can sing the refrain. he's learned about whales, India, drooping roses: i love all the new words this book introduces and the silliness of it all. this is a classic book i am willing to bet my son will STILL be able to recite from memory when he is in college

I was really looking forward to these pods based on the reviews. Starbucks is good, but I prefer bolder taste.... imagine my surprise when I ordered 2 boxes - both were expired! One expired back in 2005 for gosh sakes. I admit that Amazon agreed to credit me for cost plus part of shipping, but geez, 2 years expired!!! I'm hoping to find local San Diego area shoppe that carries pods so that I can try something different than starbucks.

Great ingredients although, chicken should have been 1st rather than chicken broth, the only thing I do not think belongs in it is Canola oil. Canola or rapeseed is not something a dog would ever find in nature and if it did find rapeseed in nature and eat it, it would poison them. Today's Food industries have convinced the masses that Canola oil is a safe and even better oil than olive or virgin coconut, facts though say otherwise. Until the late 70's it was poisonous until they figured out a way to fix that. I still like it but it could be better.

Can't do sugar. Have tried scores of SF Syrups. NONE of them can touch the excellence of this product. Thick, delicious. Perfect. 3 ingredients: Water, Maltitol, Natural Maple Flavor. PERIOD. No chemicals. No garbage. Have numerous friends & family members hooked on this stuff. My husband & son, who do NOT like "sugar free" prefer this over major label regular syrup. I use this as my SWEETENER in baking: cheesecakes, white brownies, muffins, pumpkin pies, etc... Unbelievably delicious...Can you tell I like it? :)

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 ingredients although, chicken should have been 1st rather than chicken broth, the only thing I do not think belongs in it is Canola oil. Canola or rapeseed is not something a dog would ever fi

I do not think belongs in it is Canola oil. Canola or rapeseed is not something a dog would ever find in nature and if it did find rapeseed in nature and eat it, it would poison them. Today is Food industries have convinced the masses that Canola oil is a safe and even better oil than olive or virgin coconut, facts though say otherwise. Until the late 70 is it was poisonous until they figured out a way to fix that. I still like it but it could be better.

In [19]:

```
#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 witty little book makes my son laugh at loud. i recite it in the car as we're driving along and he always can sing the refrain. he's learned about whales, India, drooping roses: i love all the new words this book introduces and the silliness of it all. this is a classic book i am willing to bet my son will STILL be able to recite from memory when he is in college

In [20]:

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

Great ingredients although chicken should have been 1st rather than chicken broth the only thing I do not think belongs in it is Canola oil Canola or rapeseed is not something a dog would ever find in nature and if it did find rapeseed in nature and eat it it would poison them Today is Food industries have convinced the masses that Canola oil is a safe and even better oil than olive or virgin coconut facts though say otherwise Until the late 70 is it was poisonous until they figured out a way to fix that I still like it but it could be better

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%|██| 364171/364171  
[03:06<00:00, 1948.15it/s]

In [23]:

```
data = preprocessed_reviews[:100000]
print(len(data))
```

100000

## [4] Featurization

## [4.1] TF-IDF

In [24]:

```
# TF-IDF
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10, use_idf = True)
tf_idf_vect.fit(data)
tf_idf_data = tf_idf_vect.fit_transform(data)
```

## [5] Assignment 11: Truncated SVD

1. **Apply Truncated-SVD on only this feature set:**

- **SET 2:** Review text, preprocessed one converted into vectors using (TFIDF)
- **Procedure:**
  - Take top 2000 or 3000 features from tf-idf vectorizers using `idf_score`.
  - You need to calculate the co-occurrence matrix with the selected features (Note:  $X.X^T$  doesn't give the co-occurrence matrix, it returns the covariance matrix, check these blogs [blog-1](#), [blog-2](#) for more information)
  - You should choose the `n_components` in truncated svd, with maximum explained variance. Please search on how to choose that and implement them. (hint: plot of cumulative explained variance ratio)
  - After you are done with the truncated svd, you can apply K-Means clustering and choose the best number of clusters based on elbow method.
  - Print out wordclouds for each cluster, similar to that in previous assignment.
  - You need to write a function that takes a word and returns the most similar words using cosine similarity between the vectors(vector: a row in the matrix after truncatedSVD)

## Truncated-SVD

### [5.1] Taking top features from TFIDF, SET 2

In [25]:

```
# Top 3000 idfs
idf = tf_idf_vect.idf_
top_3000_feature_idfs = np.argsort(idf)[-3000:]
```

In [26]:

```
# Top 3000 features
top_3000_features = np.take(tf_idf_vect.get_feature_names(), top_3000_feature_ids)
data_3000_features = tf_idf_data[:, top_3000_feature_ids]
```

## [5.2] Calculation of Co-occurrence matrix

In [96]:

```
# Ref: https://stackoverflow.com/questions/42814452/co-occurrence-matrix-from-list-of-words-in-pyth
hon
import itertools
from scipy.sparse import csr_matrix

def create_co_occurrences_matrix(allowed_words, documents):
    #print(f"allowed_words:\n{allowed_words}")
    #print(f"documents:\n{documents}")
    word_to_id = dict(zip(allowed_words, range(len(allowed_words))))
    #print(f"word_to_id:\n{word_to_id}")
    documents_as_ids = [np.sort([word_to_id[w] for w in allowed_words if doc.count(w) > 0 for i in
range(doc.count(w))]).astype('uint32') for doc in documents]
    #print(f"documents_as_ids:\n{documents_as_ids}")
    row_ind, col_ind = zip(*itertools.chain(*[(i, w) for w in doc for i, doc in
enumerate(documents_as_ids)]))
    #print(f"row_ind:\n{row_ind}")
    #print(f"col_ind:\n{col_ind}")
    data = np.ones(len(row_ind), dtype='uint32') # use unsigned int for better memory utilization
    #print(f"data:\n{data}")
    max_word_id = max(itertools.chain(*documents_as_ids)) + 1
    #print(f"max_word_id:\n{max_word_id}")
    docs_words_matrix = csr_matrix((data, (row_ind, col_ind)), shape=(len(documents_as_ids), max_wor
rd_id)) # efficient arithmetic operations with CSR * CSR
    #print(f"docs_words_matrix:\n{docs_words_matrix}")
    words_cooc_matrix = docs_words_matrix.T * docs_words_matrix # multiplying docs_words_matrix wi
th its transpose matrix would generate the co-occurrences matrix
    #print(f"words_cooc_matrix:\n{words_cooc_matrix}")
    words_cooc_matrix.setdiag(0)
    #print(f"words_cooc_matrix:\n{words_cooc_matrix.todense()}")
    return words_cooc_matrix, word_to_id
```

In [97]:

```
words = ["ABC", "PQR", "DEF"]
doc = ["ABC DEF IJK PQR", "PQR KLM OPQ", "LMN PQR XYZ ABC DEF PQR ABC"]
cooc_matrix, ids = create_co_occurrences_matrix(words, doc)
```

In [98]:

```
print(f"words_cooc_matrix:\n{cooc_matrix.todense()}")
```

```
words_cooc_matrix:
[[0 5 3]
 [5 0 3]
 [3 3 0]]
```

In [99]:

```
co_occurrence_matrix, word_to_id = create_co_occurrences_matrix(top_3000_features, data)
```

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

In [100]:

```
from sklearn.decomposition import TruncatedSVD
# components
components = range(50, 2951, 50)
```

```

components = range(0, 3000, 500)
explained_variance_ratio = []
for c in components:
    svd = TruncatedSVD(n_components=c)
    svd.fit(co_occurrence_matrix)
    explained_variance_ratio.append(svd.explained_variance_ratio_.sum())

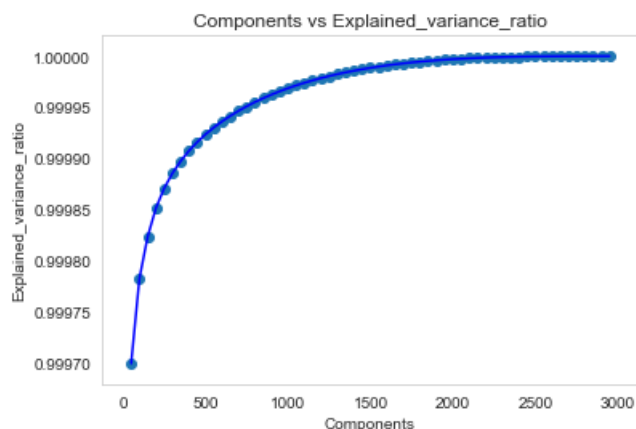
```

In [101]:

```

# plot components vs Explained_variance_ratio
sns.set_style("whitegrid", {'axes.grid' : False})
plt.plot(components, explained_variance_ratio, "b-")
plt.scatter(components, explained_variance_ratio)
plt.xlabel("Components")
plt.ylabel("Explained_variance_ratio")
plt.title("Components vs Explained_variance_ratio")
plt.show()

```



In [102]:

```

n_components = 2500
data_kmeans = co_occurrence_matrix[:, :n_components]
print(data_kmeans.shape)

```

(3000, 2500)

In [103]:

```

svd = TruncatedSVD(n_components=2500)
svd.fit(co_occurrence_matrix)
print(svd.explained_variance_ratio_)
print(svd.explained_variance_ratio_.sum())

```

```

[8.10938976e-01 8.09255216e-02 5.35688691e-02 ... 1.90800929e-09
 1.88514020e-09 1.87144814e-09]
0.99999969626152

```

## [5.4] Applying k-means clustering

In [104]:

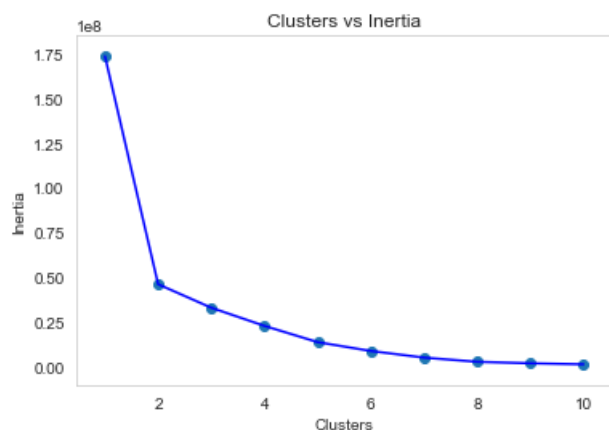
```

from sklearn.cluster import KMeans

# Clusters
K_clusters = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
# apply kmeans
inertia = []
for k in K_clusters:
    clf = KMeans(n_clusters = k, random_state = 42)
    clf.fit(data_kmeans)
    clf.predict(data_kmeans)
    inertia.append(clf.inertia_)

```

```
# plot clusters vs inertia
sns.set_style("whitegrid",{ 'axes.grid' : False})
plt.plot(K_clusters,inertia,"b-")
plt.scatter(K_clusters,inertia)
plt.xlabel("Clusters")
plt.ylabel("Inertia")
plt.title("Clusters vs Inertia")
plt.show()
```



In [105]:

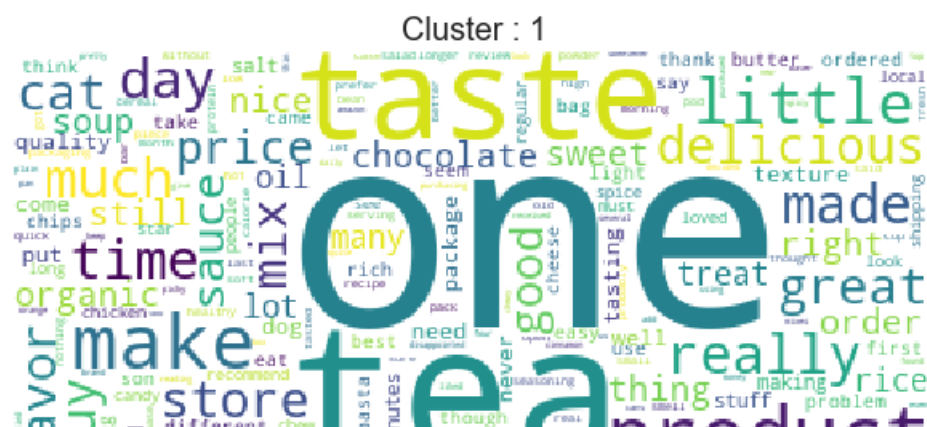
```
# best K
k = 10

# Get cluster labels
clf = KMeans(n_clusters = k, random_state = 42)
clf.fit(data_kmeans)
clf.predict(data_kmeans)
cluster_centers = clf.cluster_centers_
cluster_labels = clf.labels_.tolist()
```

## [5.5] Wordclouds of clusters obtained in the above section

In [106]:

```
from wordcloud import WordCloud
# WordCloud
for i in range(0,k):
    index = [j for j, val in enumerate(cluster_labels) if val==i]
    features = list(np.array(top_3000_features)[index])
    words = " ".join(feature for feature in features)
    WC = WordCloud(background_color='white').generate(words)
    # plot the WordCloud image
    plt.figure(figsize = (8, 8), facecolor = None)
    plt.imshow(WC)
    plt.axis("off")
    plt.tight_layout(pad = 0)
    plt.title("Cluster : "+str(i+1), fontdict={'fontsize' : 20, 'fontweight' : 20})
    plt.show()
```





Cluster : 2

sur

Cluster : 3

redi

Cluster : 4

SW

Cluster : 5

ca1

cat

Cluster : 6

char chased  
shin comb

Cluster : 7

org

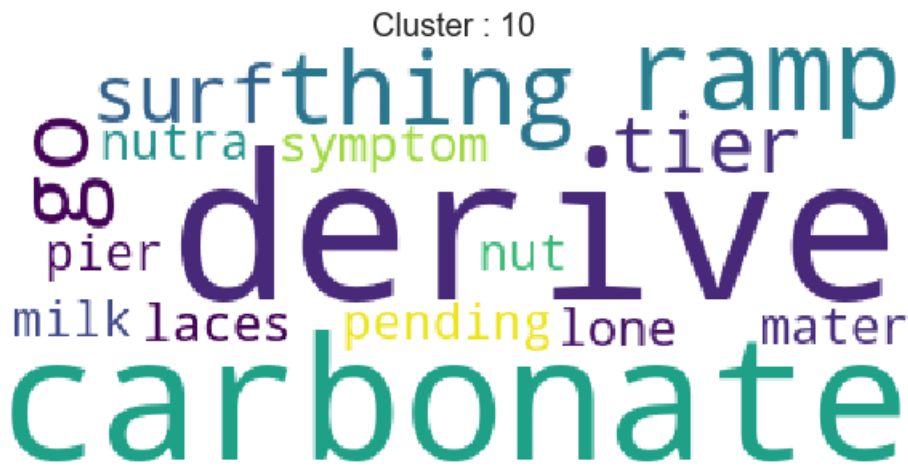
Cluster : 8

natura

Cluster : 9

com

# sap nv



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

In [109]:

```
# Function - To get most similar words
from sklearn.metrics.pairwise import cosine_similarity

def getMostSimilarWords(data, features, word, similarity = 0.9):
    index = list(features).index(word)
    vect = data[index]
    #similar_words = []
    cs = cosine_similarity(data, vect)
    #cs = cs[cs>=similarity]
    indices = [i for i, x in enumerate(cs) if x >= similarity]
    similar_words = list(np.array(features)[indices])

    return similar_words

word = top_3000_features[100]
print(word)
similar_words_lst = getMostSimilarWords(data_kmeans, top_3000_features, word, 0.75)
print(similar_words_lst)
```

```
degree water
['day though', 'degree water', 'came mind', 'smaller not', 'seems taste', 'really expect', 'thing
tea', 'sure please', 'sure drinking', 'matcha powder', 'like power', 'like syrup', 'least cup',
'loved think', 'longer brew', 'ordered mistake']
```

## [6] Conclusions

1. In TruncatedSVD as increasing the `n_components` we are getting more `explaine_variance_ratio` but after some components the `value of explaine_variance_ratio` is not increasing. It remains constant.
2. In kmeans cluster, most clusters have very few words and very few clusters have more wor

```
ds(Ex:cluster5).
```

3. Below are the Obtained values:

Optimal Explaine\_variance\_ration = 0.99

Components = 2500

Clusters = 10



In [ ]: