

Amazon Fine Food Reviews Analysis

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

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

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

Number of reviews: 568,454

Number of users: 256,059

Number of products: 74,258

Timespan: Oct 1999 - Oct 2012

Number of Attributes/Columns in data: 10

Attribute Information:

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

Objective:

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

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

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

[1]. Reading Data

[1.1] Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

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

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

In [166]:

```
%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 [167]:

```
# 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 L
# IMIT 500000""", con)
# for tsne assignment you can take 5k data points

filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3 LIM
IT 100000""", con)

# Give reviews with Score>3 a positive rating(1), and reviews with a score<3 a n
egative rating(0).
def partition(x):
    if x < 3:
        return 0
    return 1

#changing reviews with score less than 3 to be positive and vice-versa
actualScore = filtered_data['Score']
positiveNegative = actualScore.map(partition)
filtered_data['Score'] = positiveNegative
print("Number of data points in our data", filtered_data.shape)
filtered_data.head(3)
```

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

Out[167]:

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

In [168]:

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

In [169]:

```
print(display.shape)
display.head()
```

(80668, 7)

Out[169]:

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

In [170]:

```
display[display['UserId']=='AZY10LLTJ71NX']
```

Out[170]:

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

In [171]:

```
display['COUNT(*)'].sum()
```

Out[171]:

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 [172]:

```
display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID
""", con)
display.head()
```

Out[172]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	Helpfulness
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan	2	

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

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

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

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

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

In [173]:

```
#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 [174]:

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

Out[174]:

(87775, 10)

In [175]:

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

Out[175]:

87.775

Observation:- It was also seen that in two rows given below the value of HelpfulnessNumerator is greater than HelpfulnessDenominator which is not practically possible hence these two rows too are removed from calculations

In [176]:

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

display.head()
```

Out[176]:

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

In [177]:

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

In [178]:

```
#Before starting the next phase of preprocessing lets see the number of entries left
print(final.shape)

#How many positive and negative reviews are present in our dataset?
final['Score'].value_counts()
```

(87773, 10)

Out[178]:

```
1    73592
0    14181
Name: Score, dtype: int64
```

[3] Preprocessing

[3.1]. Preprocessing Review Text

Now that we have finished deduplication our data requires some preprocessing before we go on further with analysis and making the prediction model.

Hence in the Preprocessing phase we do the following in the order below:-

1. Begin by removing the html tags
2. Remove any punctuations or limited set of special characters like , or . or # etc.
3. Check if the word is made up of english letters and is not alpha-numeric
4. Check to see if the length of the word is greater than 2 (as it was researched that there is no adjective in 2-letters)
5. Convert the word to lowercase
6. Remove Stopwords
7. Finally Snowball Stemming the word (it was observed to be better than Porter Stemming)

After which we collect the words used to describe positive and negative reviews

In [179]:

```
# printing some random reviews
sent_0 = final['Text'].values[0]
print(sent_0)
print("="*50)

sent_1000 = final['Text'].values[1000]
print(sent_1000)
print("="*50)

sent_1500 = final['Text'].values[1500]
print(sent_1500)
print("="*50)

sent_4900 = final['Text'].values[4900]
print(sent_4900)
print("="*50)
```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its very hard to find any chicken products made in the USA but they are out there, but this one isnt. Its too bad too because its a good product but I wont take any chances till they know what is going on with the china imports.

=====

The Candy Blocks were a nice visual for the Lego Birthday party but the candy has little taste to it. Very little of the 2 lbs that I bought were eaten and I threw the rest away. I would not buy the candy again.

=====

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

=====

My dog LOVES these treats. They tend to have a very strong fish oil smell. So if you are afraid of the fishy smell, don't get it. But I think my dog likes it because of the smell. These treats are really small in size. They are great for training. You can give your dog several of these without worrying about him over eating. Amazon's price was much more reasonable than any other retailer. You can buy a 1 pound bag on Amazon for almost the same price as a 6 ounce bag at other retailers. It's definitely worth it to buy a big bag if your dog eats them a lot.

In [180]:

```
# 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_1500 = re.sub(r"http\S+", "", sent_1500)
sent_4900 = re.sub(r"http\S+", "", sent_4900)

print(sent_0)
```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its very hard to find any chicken products made in the USA but they are out there, but this one isnt. Its too bad too because its a good product but I wont take any chances till they know what is going on with the china imports.

In [181]:

```
# https://stackoverflow.com/questions/16206380/python-beautifulsoup-how-to-remove-all-tags-from-an-element
from bs4 import BeautifulSoup

soup = BeautifulSoup(sent_0, 'lxml')
text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_1000, 'lxml')
text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_1500, 'lxml')
text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_4900, 'lxml')
text = soup.get_text()
print(text)
```

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its very hard to find any chicken products made in the USA but they are out there, but this one isnt. Its too bad too because its a good product but I wont take any chances till they know what is going on with the china imports.

=====

The Candy Blocks were a nice visual for the Lego Birthday party but the candy has little taste to it. Very little of the 2 lbs that I bought were eaten and I threw the rest away. I would not buy the candy again.

=====

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

=====

My dog LOVES these treats. They tend to have a very strong fish oil smell. So if you are afraid of the fishy smell, don't get it. But I think my dog likes it because of the smell. These treats are really small in size. They are great for training. You can give your dog several of these without worrying about him over eating. Amazon's price was much more reasonable than any other retailer. You can buy a 1 pound bag on Amazon for almost the same price as a 6 ounce bag at other retailers. It's definitely worth it to buy a big bag if your dog eats them a lot.

In [182]:

```
# 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 [183]:

```
sent_1500 = decontracted(sent_1500)
print(sent_1500)
print("="*50)
```

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

In [184]:

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

My dogs loves this chicken but its a product from China, so we wont be buying it anymore. Its very hard to find any chicken products made in the USA but they are out there, but this one isnt. Its too bad too because its a good product but I wont take any chances till they know what is going on with the china imports.

In [185]:

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

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

In [186]:

```
# 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 revmoved in the 1st s
tep

stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ou
rselves', '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', 'itse
lf', 'they', 'them', 'their', \
               'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'tha
t', "that'll", 'these', 'those', \
               'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'ha
s', 'had', 'having', 'do', 'does', \
               'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'becaus
e', 'as', 'until', 'while', 'of', \
               'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 't
hrough', 'during', 'before', 'after', \
               'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'of
f', 'over', 'under', 'again', 'further', \
               'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'al
l', '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'v
e", 'now', 'd', 'll', 'm', 'o', 're', \
               've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "d
idn'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', "should
n't", 'wasn', "wasn't", 'weren', "weren't", \
               'won', "won't", 'wouldn', "wouldn't"])
```

In [187]:

```
# Combining all the above stundents
from tqdm import tqdm
preprocessed_reviews = []
# tqdm is for printing the status bar
for sentence in tqdm(final['Text'].values):
    sentence = re.sub(r"http\S+", "", sentence)
    sentence = BeautifulSoup(sentence, 'lxml').get_text()
    sentence = decontracted(sentence)
    sentence = re.sub("\S*\d\S*", "", sentence).strip()
    sentence = re.sub('[^A-Za-z]+', '', sentence)
    # https://gist.github.com/sebleier/554280
    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not in
stopwords)
    preprocessed_reviews.append(sentence.strip())
```

100%|██████████| 87773/87773 [00:33<00:00, 2647.52it/s]

In [188]:

```
preprocessed_reviews[1500]
```

Out[188]:

```
'way hot blood took bite jig lol'
```

In [189]:

```
len(preprocessed_reviews)
```

Out[189]:

```
87773
```

In [190]:

```
def do_pickling(filename,data):
    with open(filename, "wb") as f:
        pickle.dump(data,f)
```

In [191]:

```
do_pickling('final.pickle',preprocessed_reviews)
```

[3.2] Preprocessing Review Summary

In [6]:

```
## Similarly you can do preprocessing for review summary also.
```

[4] Featurization

[4.1] BAG OF WORDS

In [25]:

```
#BoW
count_vect = CountVectorizer() #in scikit-learn
count_vect.fit(preprocessed_reviews)
print("some feature names ", count_vect.get_feature_names()[:10])
print('='*50)

final_counts = count_vect.transform(preprocessed_reviews)
print("the type of count vectorizer ",type(final_counts))
print("the shape of out text BOW vectorizer ",final_counts.get_shape())
print("the number of unique words ", final_counts.get_shape()[1])
```

```
some feature names  ['aa', 'aahhhs', 'aback', 'abandon', 'abates',
'abbott', 'abby', 'abdominal', 'abiding', 'ability']
=====
```

```
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (4986, 12997)
the number of unique words  12997
```

[4.2] Bi-Grams and n-Grams.

In [26]:

```
#bi-gram, tri-gram and n-gram

#removing stop words like "not" should be avoided before building n-grams
# count_vect = CountVectorizer(ngram_range=(1,2))
# please do read the CountVectorizer documentation http://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html

# you can choose these numebers min_df=10, max_features=5000, of your choice
count_vect = CountVectorizer(ngram_range=(1,2), min_df=10, max_features=5000)
final_bigram_counts = count_vect.fit_transform(preprocessed_reviews)
print("the type of count vectorizer ",type(final_bigram_counts))
print("the shape of out text BOW vectorizer ",final_bigram_counts.get_shape())
print("the number of unique words including both unigrams and bigrams ", final_b
igram_counts.get_shape()[1])
```

```
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer (4986, 3144)
the number of unique words including both unigrams and bigrams 3144
```

[4.3] TF-IDF

In [27]:

```
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
tf_idf_vect.fit(preprocessed_reviews)
print("some sample features(unique words in the corpus)",tf_idf_vect.get_feature
_names()[0:10])
print('='*50)

final_tf_idf = tf_idf_vect.transform(preprocessed_reviews)
print("the type of count vectorizer ",type(final_tf_idf))
print("the shape of out text TFIDF vectorizer ",final_tf_idf.get_shape())
print("the number of unique words including both unigrams and bigrams ", final_t
f_idf.get_shape()[1])
```

```
some sample features(unique words in the corpus) ['ability', 'able',
'able find', 'able get', 'absolute', 'absolutely', 'absolutely delic
ious', 'absolutely love', 'absolutely no', 'according']
```

```
=====
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer (4986, 3144)
the number of unique words including both unigrams and bigrams 3144
```

In [192]:

```
tf_idf_vect = TfidfVectorizer(ngram_range=(1,1),min_df=10,max_df =2000)
tf_idf_train = tf_idf_vect.fit_transform(preprocessed_reviews)
```

In [193]:

```
do_pickling('tfidf.pickle',tf_idf_train)
do_pickling('tfidfvect.pickle',tf_idf_vect)
```

[4.4] Word2Vec

In [31]:

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

In [33]:

```
# Using Google News Word2Vectors

# in this project we are using a pretrained model by google
# its 3.3G file, once you load this into your memory
# it occupies ~9Gb, so please do this step only if you have >12G of ram
# we will provide a pickle file wich contains a dict ,
# and it contains all our courpus words as keys and model[word] as values
# To use this code-snippet, download "GoogleNews-vectors-negative300.bin"
# from https://drive.google.com/file/d/0B7XkCwpI5KDYNlNUTTlSS21pQmM/edit
# it's 1.9GB in size.

# http://kavita-ganesan.com/gensim-word2vec-tutorial-starter-code/#.W17SRFAzZPY
# you can comment this whole cell
# or change these variable according to your need

is_your_ram_gt_16g=False
want_to_use_google_w2v = False
want_to_train_w2v = True

if want_to_train_w2v:
    # min_count = 5 considers only words that occured atleast 5 times
    w2v_model=Word2Vec(list_of_sentence,min_count=5,size=50, workers=4)
    print(w2v_model.wv.most_similar('great'))
    print('='*50)
    print(w2v_model.wv.most_similar('worst'))

elif want_to_use_google_w2v and is_your_ram_gt_16g:
    if os.path.isfile('GoogleNews-vectors-negative300.bin'):
        w2v_model=KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative
300.bin', binary=True)
        print(w2v_model.wv.most_similar('great'))
        print(w2v_model.wv.most_similar('worst'))
    else:
        print("you don't have gogole's word2vec file, keep want_to_train_w2v = T
rue, to train your own w2v ")

[('fantastic', 0.8510183691978455), ('excellent', 0.835196912288665
8), ('good', 0.8275485634803772), ('terrific', 0.8223178386688232),
('awesome', 0.8122442960739136), ('wonderful', 0.7793042063713074),
('perfect', 0.7753506898880005), ('nice', 0.7478526830673218), ('fab
ulous', 0.7191284894943237), ('amazing', 0.7103496789932251)]

=====
[('greatest', 0.8144274950027466), ('tastiest', 0.7520067095756531),
('best', 0.7180012464523315), ('nastiest', 0.6530043482780457), ('di
sgusting', 0.630192756652832), ('smoothest', 0.6221706867218018),
('coolest', 0.6177399158477783), ('closest', 0.6112560033798218),
('surpass', 0.5932056903839111), ('horrible', 0.5899286866188049)]
```


In [34]:

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

```
number of words that occurred minimum 5 times 17386
sample words ['dogs', 'loves', 'chicken', 'product', 'china', 'won
t', 'buying', 'anymore', 'hard', 'find', 'products', 'made', 'usa',
'one', 'isnt', 'bad', 'good', 'take', 'chances', 'till', 'know', 'go
ing', 'imports', 'love', 'saw', 'pet', 'store', 'tag', 'attached',
'regarding', 'satisfied', 'safe', 'infestation', 'literally', 'every
where', 'flying', 'around', 'kitchen', 'bought', 'hoping', 'least',
'get', 'rid', 'weeks', 'fly', 'stuck', 'squishing', 'buggers', 'succ
ess', 'rate']
```

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

[4.4.1.1] Avg W2v

In [38]:

```
# average Word2Vec
# compute average word2vec for each review.
sent_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(list_of_sentence): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need
    # to change this to 300 if you use google's w2v
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    sent_vectors.append(sent_vec)
print(len(sent_vectors))
print(len(sent_vectors[0]))
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 4986/4986 [00:03<00:00, 1330.47it/s]
```

```
4986
50
```

[4.4.1.2] TFIDF weighted W2v

In [35]:

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

[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](https://medium.com/data-science-group-iitr/word-embedding-2d05d270b285), (<https://medium.com/data-science-group-iitr/word-embedding-2d05d270b285>) [blog-2](https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/) (<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/>) for more information)
 - You should choose the `n_components` in truncated svd, with maximum explained variance. Please search on how to choose that and implement them. (hint: plot of cumulative explained variance ratio)
 - After you are done with the truncated svd, you can apply K-Means clustering and choose the best number of clusters based on elbow method.
 - Print out wordclouds for each cluster, similar to that in previous assignment.
 - You need to write a function that takes a word and returns the most similar words using cosine similarity between the vectors (vector: a row in the matrix after truncatedSVD)

Truncated-SVD

In [194]:

```
import sqlite3
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import pickle

from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.decomposition import TruncatedSVD
from sklearn.utils.extmath import randomized_svd
from sklearn.metrics.pairwise import cosine_similarity
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from collections import defaultdict
from wordcloud import WordCloud
from scipy.spatial.distance import pdist, squareform
```

In [195]:

```
# function to load the pickle data
def loadPickleData(filename):
    pickle_off = open(filename, "rb")
    final = pickle.load(pickle_off)
    return final
```

In [196]:

```
def plot_find_best_k(train):
    inertia_list = []
    for i in range(5,50,5):
        kmeans = KMeans(n_clusters=i, random_state=0).fit(train)
        inertia_list.append(kmeans.inertia_)
    fig = plt.figure()
    ax1 = fig.add_subplot(1,1,1)
    ax1.plot(inertia_list)
```

In [197]:

```
def plot(explained_variance):
    fig = plt.figure(figsize=(10,10))
    ax1 = fig.add_subplot(1,1,1)
    ax1.plot(explained_variance.cumsum(), 'k--')
```

In [198]:

```
# select the no of components giving the requiried variance
def select_n_components(var_ratio, goal_var: float) -> int:
    # Set initial variance explained so far
    total_variance = 0.0

    # Set initial number of features
    n_components = 0

    # For the explained variance of each feature:
    for explained_variance in var_ratio:

        # Add the explained variance to the total
        total_variance += explained_variance

        # Add one to the number of components
        n_components += 1

        # If we reach our goal level of explained variance
        if total_variance >= goal_var:
            # End the loop
            break

    # Return the number of components
    return n_components
```

In [199]:

```
def getTextIndices(labels):
    indices_dict = defaultdict(lambda: [])
    for i,x in enumerate(labels):
        indices_dict[x].append(i)
    return indices_dict
```

In [200]:

```
def do_pickling(filename,data):
    with open(filename, "wb") as f:
        pickle.dump(data,f)
```

In [201]:

```
def plotWordCloudForClusters(text_dict):
    for key,value in text_dict.items():
        print("For cluster ",key)
        draw_wordcloud(value)
```

In [202]:

```
def draw_wordcloud(train):
    wordcloud = WordCloud(        background_color='black',
                                width=1600,
                                height=800,
                                ).generate(train)

    fig = plt.figure(figsize=(30,20))
    plt.imshow(wordcloud)
    plt.axis('off')
    plt.tight_layout(pad=0)
    plt.show()
```

In [203]:

```
def flatten_the_list(text_dict):
    for key,value in text_dict.items():
        text = ''
        for i in value:
            text = text + ' ' + i
        text_dict[key]=text
    return text_dict
```

[5.1] Taking top features from TFIDF, SET 2

In [204]:

```
# load the pickle data
final = loadPickleData('final.pickle')
tfidf = loadPickleData('tfidf.pickle')
tf_idf_vect = loadPickleData('tfidfvect.pickle')
```

In [206]:

```
# getting indices of top features having higher idf values
indices = np.argsort(-tf_idf_vect.idf_)[:3000]
```

In [207]:

```
# mapping indices to feature names
feature_names = tf_idf_vect.get_feature_names()
feature_list = []
for x in indices:
    feature_list.append(feature_names[x])
print(feature_list[:10])
```

```
['objective', 'stared', 'cheer', 'malamute', 'uplifting', 'fierce',
'fiesta', 'painless', 'fertilizers', 'painfully']
```

In [218]:

```
words = feature_list
```

[5.2] Calculation of Co-occurrence matrix

In [220]:

```
reviews = final
co_occurence_matrix = np.zeros((len(words),len(words)))
window = 5
```

In [221]:

```
len(reviews)
```

Out[221]:

87773

In [222]:

```
print(co_occurence_matrix.shape)
```

(3000, 3000)

In [223]:

```
for review in reviews:
    text = review.split()
    for i,word in enumerate(text):
        if(word in words):
            # go backwards till window length
            j=i-1
            while(j>=0 and j>=i-window):
                word_col = text[j]
                if(word_col in words):
                    l = words.index(word)
                    m = words.index(word_col)
                    co_occurence_matrix[l][m]+=1
                j=j-1
            #go forward till the window length
            k=i+1
            while(k<len(text) and k<=i+window):
                word_col = text[k]
                if(word_col in words):
                    l = words.index(word)
                    q = words.index(word_col)
                    co_occurence_matrix[l][q]+=1
                k=k+1

for i,word in enumerate(words):
    co_occurence_matrix[i][i]+=1
```

In [224]:

```
do_pickling('co-occurencematrix.pickle',co_occurence_matrix)
```

In [225]:

```
co_occurence_matrix = loadPickleData('co-occurencematrix.pickle')
```

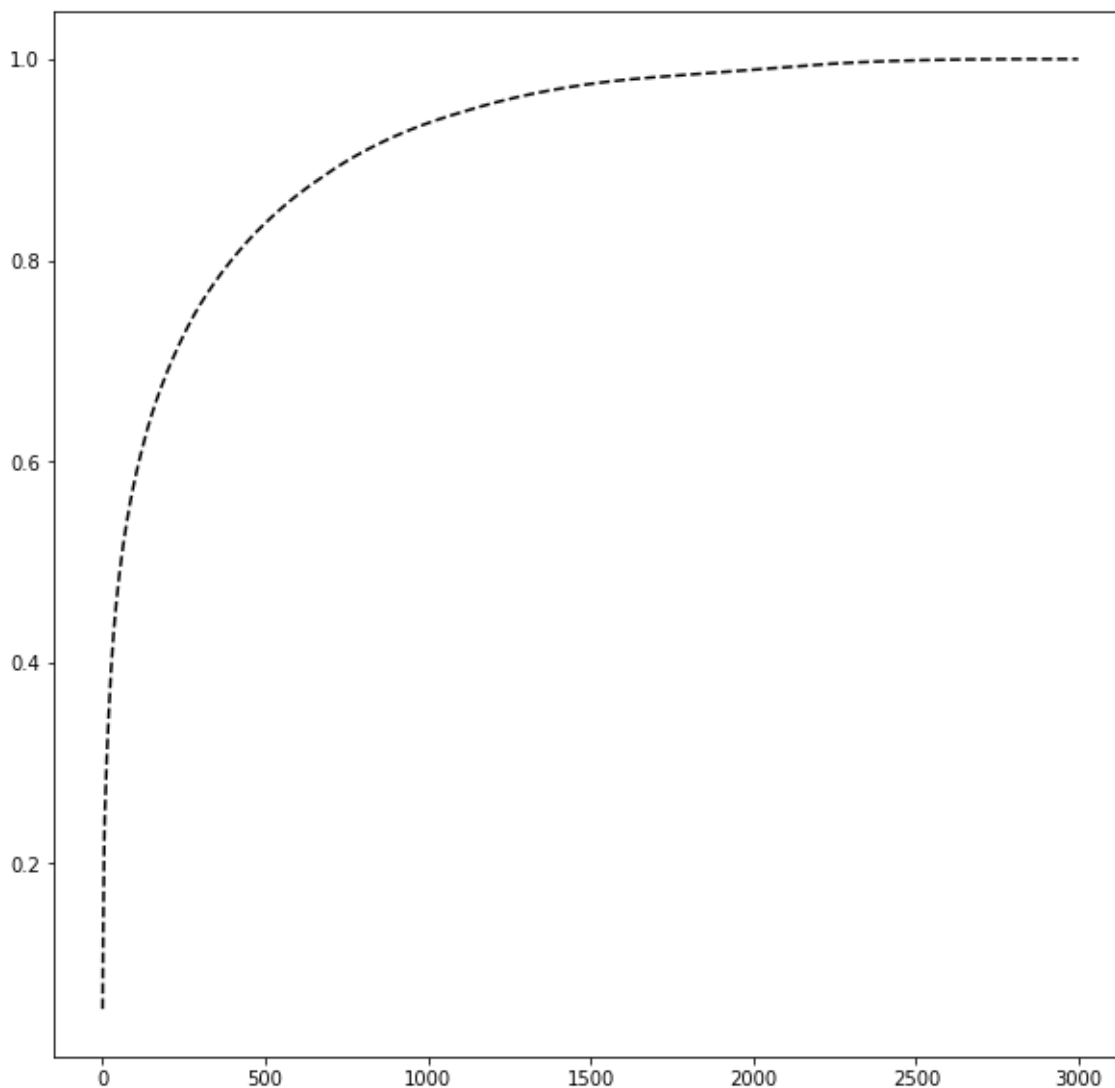
[5.3] Finding optimal value for number of components (n) to be retained.

In [226]:

```
# Please write all the code with proper documentation  
# applying truncated svd  
trun_svd = TruncatedSVD(n_components=co_occurence_matrix.shape[1]-1, algorithm=  
'randomized', n_iter=5, random_state=None, tol=0.0)  
trun_svd.fit(co_occurence_matrix)  
explained_variance = trun_svd.explained_variance_ratio_
```

In [227]:

```
plot(explained_variance)
```



In [228]:

```
# finding the number of components which gives 99.99 % variance
n_components = select_n_components(explained_variance, 0.99)
print(n_components)
```

2019

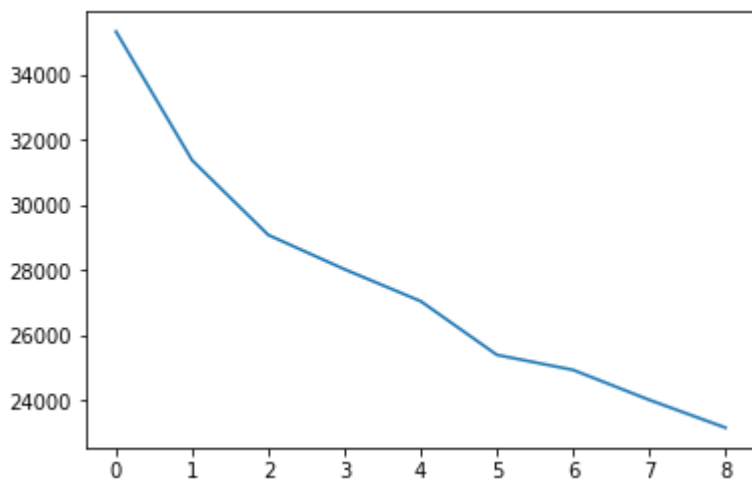
In [229]:

```
# performing truncated svd with n_components giving maximum variance
trun_svd = TruncatedSVD(n_components=n_components, algorithm='randomized', n_iter=5, random_state=None, tol=0.0)
X = trun_svd.fit_transform(co_occurrence_matrix)
```

[5.4] Applying k-means clustering

In [230]:

```
# find the best number of clusters using the elbow method
plot_find_best_k(X)
```



In [231]:

```
# applying kmeans clustering with the best k
kmeans = KMeans(n_clusters=5, init='k-means++', random_state=0).fit(X)
```

[5.5] Wordclouds of clusters obtained in the above section

In [232]:

```
def getText(indices_dict, data):
    text_dict = defaultdict(lambda: [])
    for key, value in indices_dict.items():
        for i in value:
            x = data[i]
            text_dict[key].append(x)
    text_dict = flatten_the_list(text_dict)
    return text_dict
```

In [233]:

```
# given a word and the cosine similarity matrix it returns the words having maximum cosine similarity
def findSimilarWords(word,similarity):
    index = word_to_index[word]
    indices = np.argsort(similarity[index])
    similar_words = []
    for x in indices:
        similar_words.append(index_to_word[x])
    return similar_words
```

In [234]:

```
# mapping from word to index
word_to_index = dict()
for i,x in enumerate(words):
    word_to_index[x]=i
```

In [235]:

```
# inverted mapping from index to word
index_to_word = {v :k for k,v in word_to_index.items() }
```

In [236]:

```
labels = kmeans.labels_
indices_dict = getTextIndices(labels)
```

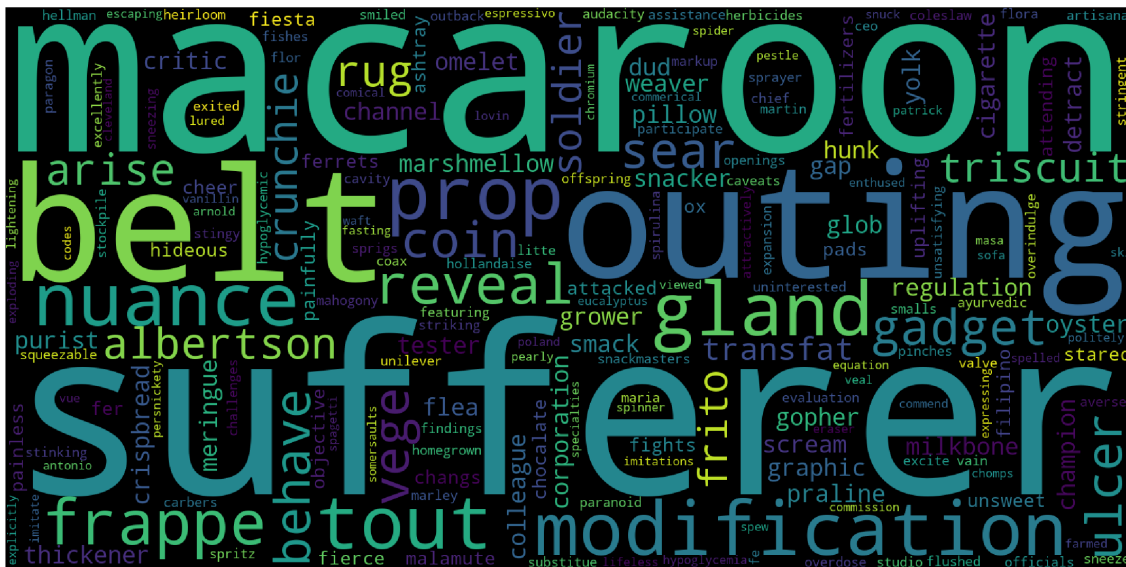
In [237]:

```
text_dict = getText(indices_dict,words)
```


In [238]:

```
plotWordCloudForClusters(text_dict)
```

For cluster 2



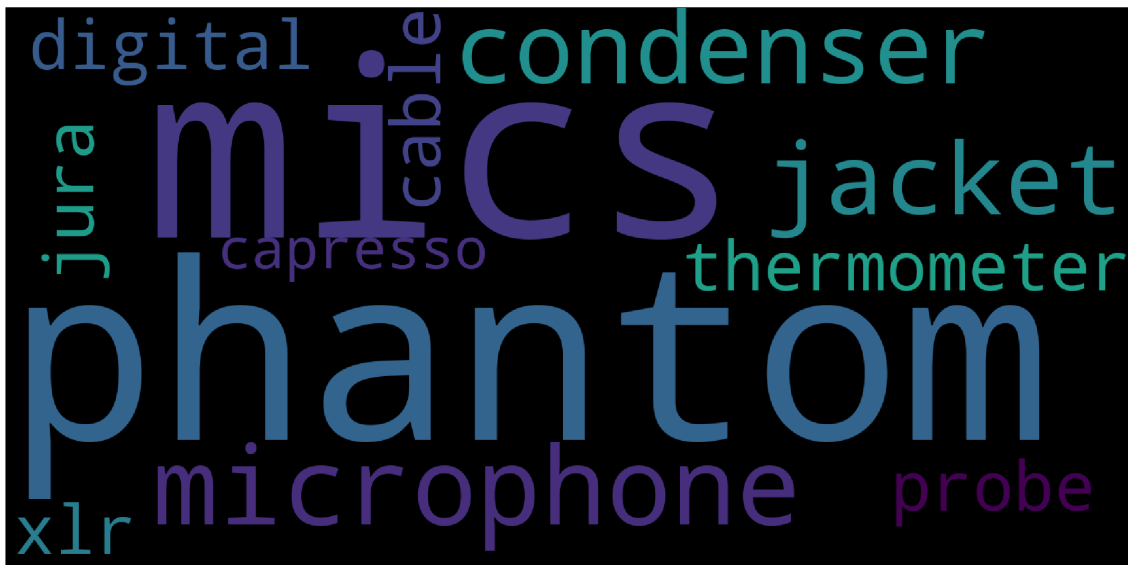
For cluster 4

proteinate

For cluster 0

koi
ocha oi

For cluster 1



For cluster 3



In [239]:

```
print(indices_dict.keys())
```

```
dict_keys([2, 4, 0, 1, 3])
```

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

In [240]:

```
# Please write all the code with proper documentation  
'fraud' in words
```

Out[240]:

True

```
# take a word , find the cluster to which it belongs and find most similar words
s_word = 'fraud'
index = word_to_index[word]
cluster = labels[index]
# find the vectors of words in the cluster
# 1 . find the indices of words from the labels.
# 2 . get the vector for each word from X (x[index])
# check for type of labels
labels = pd.Series(labels)
c = labels[labels==cluster]
print(c.shape)
```

```
word_vect_in_cluster = []
for i in c.index:
    word_vect_in_cluster.append(X[i])

similarity = squareform(pdist(word_vect_in_cluster, 'cosine'))
similar_words = findSimilarWords(s_word,similarity)
```

->I have taken 3000 features. ->The number of components that gives 99% variance is 2019 and I have choose the number of clusters as 5. ->Out of 5 clusters ,most of the words are present in cluster 2. ->On plotting the wordcloud for most similar words to word 'fraud', I could abuse,occurence,suspicion in the word cluster

In []: