

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>)

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).

1. Import required libraries

```
In [54]: 1 import warnings
          2 warnings.filterwarnings("ignore")
```

```
In [55]: 1 %matplotlib inline
2
3 import sqlite3
4 import pandas as pd
5 import numpy as np
6 import nltk
7 import string
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10 from sklearn.feature_extraction.text import TfidfTransformer
11 from sklearn.feature_extraction.text import TfidfVectorizer
12
13 from sklearn.feature_extraction.text import CountVectorizer
14 from sklearn import metrics
15 from sklearn.model_selection import train_test_split
16 import re
17 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
18 import string
19 from nltk.corpus import stopwords
20 from nltk.stem import SnowballStemmer
21 from nltk.stem.wordnet import WordNetLemmatizer
22 from gensim.models import Word2Vec
23 from gensim.models import KeyedVectors
24 import pickle
25
26 from tqdm import tqdm_notebook
27 from tqdm import tqdm
28 from bs4 import BeautifulSoup
29 import os
```

2. Read the Dataset

- Create a Connection object that represents the database. Here the data will be stored in the 'database.sqlite' file.
- Read the Dataset table using connection object where the score column != 3
- Replace the score values with 'positive' and 'negative' label.(i.e Score 1 & 2 is labeled as negative and Score 4 & 5 is labeled as positive)
- Score with value 3 is neutral.


```

In [56]: 1 # using SQLite Table to read data.
2 con = sqlite3.connect('database.sqlite')
3
4 # filtering only positive and negative reviews i.e.
5 # not taking into consideration those reviews with Score=3
6 # SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000, will give top 500000 data points
7 # you can change the number to any other number based on your computing power
8
9 # filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000""", con)
10 # for tsne assignment you can take 5k data points
11
12 filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3 LIMIT 100000""", con)
13
14 # Give reviews with Score>3 a positive rating(1), and reviews with a score<3 a negative rating(0).
15 def partition(x):
16     if x < 3:
17         return 0
18     return 1
19
20 #changing reviews with score less than 3 to be positive and vice-versa
21 actualScore = filtered_data['Score']
22 positiveNegative = actualScore.map(partition)
23 filtered_data['Score'] = positiveNegative
24 print("Number of data points in our data", filtered_data.shape)
25 filtered_data.head(3)

```

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

Out[56]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	1	1303862400	Good Quality Dog Food	I have bought several of the Vitality canned d...

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	0	1346976000	Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...
2	3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1	1	1219017600	"Delight" says it all	This is a confection that has been around a fe...

Type *Markdown* and LaTeX: α^2

```
In [57]: 1 display = pd.read_sql_query("""
2         SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
3         FROM Reviews
4         GROUP BY UserId
5         HAVING COUNT(*)>1
6         """, con)
```

```
In [58]: 1 print(display.shape)
2         display.head()
```

(80668, 7)

Out[58]:

		UserId	ProductId	ProfileName	Time	Score	Text	COUNT(*)
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-R12KPBODL2B5ZD	B007OSBE1U	Christopher P. Presta	1348617600	1	I didnt like this coffee. Instead of telling y...		2

```
In [59]: 1 display[display['UserId']=='AZY10LLTJ71NX']
```

Out[59]:

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

```
In [60]: 1 display['COUNT(*)'].sum()
```

Out[60]: 393063

4. Exploratory Data Analysis

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 [61]: 1 display= pd.read_sql_query("""
2         SELECT *
3         FROM Reviews
4         WHERE Score != 3 AND UserId="AR5J8UI46CURR"
5         ORDER BY ProductID
6         """, con)
7         display.head()
```

Out[61]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS	DELICIOUS WAFERS FIND THEM EUROPEAN WAFERS
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS	DELICIOUS WAFERS FIND THEM EUROPEAN WAFERS
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS	DELICIOUS WAFERS FIND THEM EUROPEAN WAFERS
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS	DELICIOUS WAFERS FIND THEM EUROPEAN WAFERS
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS	DELICIOUS WAFERS FIND THEM EUROPEAN WAFERS

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

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 [62]: 1 #Sorting data according to ProductId in ascending order
        2 sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True, inplace=False, kind='quicksort', na_posit
```

```
In [63]: 1 #Deduplication of entries
        2 final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"}, keep='first', inplace=False)
        3 final.shape
```

Out[63]: (87775, 10)

```
In [64]: 1 #Checking to see how much % of data still remains
        2 (final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

Out[64]: 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 [65]: 1 display= pd.read_sql_query("""
2 SELECT *
3 FROM Reviews
4 WHERE Score != 3 AND Id=44737 OR Id=64422
5 ORDER BY ProductID
6 """, con)
7
8 display.head(2)
```

Out[65]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text
0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	1	5	1224892800	Bought This for My Son at College	My son loves spaghetti so I didn't hesitate or...
1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	2	4	1212883200	Pure cocoa taste with crunchy almonds inside	It was almost a 'love at first bite' - the per...

- 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

```
In [66]: 1 final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [67]: 1 #Before starting the next phase of preprocessing Lets see the number of entries left
2 print(final.shape)
3
4 #How many positive and negative reviews are present in our dataset?
5 final['Score'].value_counts()
```

(87773, 10)

```
Out[67]: 1    73592
0    14181
Name: Score, dtype: int64
```

5. Preprocessing

[5.1]. Preprocessing Review Text and Summary

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 [68]:

```
1 # https://stackoverflow.com/a/47091490/4084039
2 import re
3
4 def decontracted(phrase):
5     # specific
6     phrase = re.sub(r"won't", "will not", phrase)
7     phrase = re.sub(r"can't", "can not", phrase)
8
9     # general
10    phrase = re.sub(r"n't", " not", phrase)
11    phrase = re.sub(r"\ 're", " are", phrase)
12    phrase = re.sub(r"\ 's", " is", phrase)
13    phrase = re.sub(r"\ 'd", " would", phrase)
14    phrase = re.sub(r"\ 'll", " will", phrase)
15    phrase = re.sub(r"\ 't", " not", phrase)
16    phrase = re.sub(r"\ 've", " have", phrase)
17    phrase = re.sub(r"\ 'm", " am", phrase)
18    return phrase
```

In [69]:

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

```

In [70]: 1 # Combining all the above students
2 from tqdm import tqdm
3 def createCleanedText(review_text, column_name):
4     sno = nltk.stem.SnowballStemmer('english') #initialising the snowball stemmer
5     preprocessed_reviews = []
6     # tqdm is for printing the status bar
7     for sentence in tqdm(review_text):
8         sentence = re.sub(r"http\S+", "", sentence) # \S=except space; + = 1 or more
9         sentence = BeautifulSoup(sentence, 'lxml').get_text() # remove links
10        sentence = decontracted(sentence) # expand short forms
11        sentence = re.sub("\S*\d\S*", "", sentence).strip() #remove words containing digits
12        sentence = re.sub('[^A-Za-z]+', ' ', sentence) # remove special char
13        # https://gist.github.com/sebleier/554280
14        sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not in stopwords)
15        preprocessed_reviews.append(sentence.strip())
16    #adding a column of CleanedText which displays the data after pre-processing of the review
17    final[column_name]=preprocessed_reviews
18

```

```

In [71]: 1 if not os.path.isfile('final.sqlite'):
2     #createCleanedText(final['Text_Summary'].values, column_name='CleanedTextSumm')
3     createCleanedText(final['Text'].values, column_name='CleanedText')
4     conn = sqlite3.connect('final.sqlite')
5     c=conn.cursor()
6     conn.text_factory = str
7     final.to_sql('Reviews', conn, schema=None, if_exists='replace', \
8                 index=True, index_label=None, chunksize=None, dtype=None)
9     conn.close()

```

100%|██████████| 87773/87773 [00:34<00:00, 2573.23it/s]

```

In [72]: 1 if os.path.isfile('final.sqlite'):
2     conn = sqlite3.connect('final.sqlite')
3     final = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3 """, conn)
4     conn.close()
5 else:
6     print("Please the above cell")

```

In [73]:

```
1 print(final.head(3))
2 final.shape
```

```

      index      Id  ProductId      UserId      ProfileName \
0  22620  24750  2734888454  A13ISQV0U9GZIC      Sandikaye
1  22621  24751  2734888454  A1C298ITT645B6  Hugh G. Pritchard
2  70677  76870  B00002N8SM  A19Q006CSFT011      Arlielle

      HelpfulnessNumerator  HelpfulnessDenominator  Score      Time \
0              1              1      0  1192060800
1              0              0      1  1195948800
2              0              0      0  1288396800

      Summary      Text \
0      made in china  My dogs loves this chicken but its a product f...
1      Dog Lover Delites  Our dogs just love them. I saw them in a pet ...
2  only one fruitfly stuck  I had an infestation of fruitflies, they were ...

      CleanedText
0  dogs loves chicken product china wont buying a...
1  dogs love saw pet store tag attached regarding...
2  infestation fruitflies literally everywhere fl...
```

Out[73]: (87773, 12)

6. Splitting data into Train and Test set

In [74]:

```
1 #sorted dataframe by time
2 '''
3 df['Time']=pd.to_datetime(final['Time'],unit='s')
4 df=df.sort_values(by="Time")
5 df.head(20)
6 '''
7 df=final.sort_values(by=['Time'])
8 #df.head(5)
```

```
In [75]: 1 #TEXT COLUMN
2 X=np.array(df['CleanedText'])
3 #SCORE COLUMN
4 y=np.array(df['Score'])
```

7. Featurization

[7.2] TF-IDF

Tf-idf stands for term frequency-inverse document frequency, and the tf-idf weight is a weight often used in information retrieval and text mining. This weight is a statistical measure used to evaluate how important a word is to a document in a collection or corpus.

- 1.TF: Term Frequency, which measures how frequently a term occurs in a document.
$$TF(t) = (\text{Number of times term } t \text{ appears in a document}) / (\text{Total number of terms in the document}).$$
- 2.IDF: Inverse Document Frequency, is a scoring of how rare the word is across documents.
$$IDF(t) = \log_e(\text{Total number of documents} / \text{Number of documents with term } t \text{ in it}).$$
- 3.The scores are a weighting where not all words are equally as important or interesting.

The scores have the effect of highlighting words that are distinct (contain useful information) in a given document. The idf of a rare term is high, whereas the idf of a frequent term is likely to be low.

```
In [76]: 1 def tfidfVector(X, max_features=None):
2     tf_idf_vect = TfidfVectorizer(min_df=5,max_features=max_features)
3     X_tfidf = tf_idf_vect.fit_transform(X)
4     print("the type of count vectorizer: ",type(X))
5     print("the shape of out text TFIDF vectorizer: ",X_tfidf.get_shape())
6     print("the number of unique words including both unigrams and bigrams: ", X_tfidf.get_shape()[1])
7     return tf_idf_vect, X_tfidf
```

```
In [77]: 1 # Tfidf vector with all features which we use for brute force implementation
        2 %time tf_idf_vect, X_tfidf=tfidfVector(X,max_features=None)
```

the type of count vectorizer: <class 'numpy.ndarray'>
the shape of out text TFIDF vectorizer: (87773, 16545)
the number of unique words including both unigrams and bigrams: 16545
CPU times: user 3.64 s, sys: 64 ms, total: 3.71 s
Wall time: 3.63 s

8. Function for object state :

- a. savetofile(): to save the current state of object for future use using pickle.
- b. openfromfile(): to load the past state of object for further use.

```
In [43]: 1 #Functions to save objects for later use and retireve it
        2 def savetofile(obj,filename):
        3     pickle.dump(obj,open(filename+".pkl","wb"))
        4 def openfromfile(filename):
        5     temp = pickle.load(open(filename+".pkl","rb"))
        6     return temp
        7
        8 savetofile(tf_idf_vect,'tf_idf_vect')
        9 savetofile(X_tfidf,'X_tfidf')
       10
       11 savetofile(X,'X')
       12 savetofile(y,'y')
```