

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

Developer : PraveenAI

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:

I wanted to generate a best n-Dimensional Datasets for Amazon Reviews. This helps the industries to get the NLTK very Handy and for developers as a reference. which wills helps me in showcasing my knowledge in Natural Processing Techniques.

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 nuetral 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""", con)
# for tsne assignment you can take 5k data points

filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3 LIMIT 70000""", 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 (70000, 10)
```

Out[2]:

	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...
1	2	B00813GRG4	A1D87F6ZCVESNK	dli pa	0	0	0	1346976000	Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...
2	3	B000LQOCH0	ABXLMWJXXAIN	Natalia Corres "Natalia Corres"	1	1	1	1219017600	"Delight" says it all	This is a confection that has been around a fe...

```
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	B007Y59HVM	Breyton	1331510400	2		Overall its just OK when considering the price...	2
1	#oc-R11D9D7SHXUJB9	B005HG9ET0	Louis E. Emory "hoppy"	1342396800	5		My wife has recurring extreme muscle spasms. u...	3
2	#oc-R11DNUJ2NBKQ23Z	B007Y59HVM	Kim Cieszykowski	1348531200	1		This coffee is horrible and unfortunately not ...	2
3	#oc-R11O5JSZVQE25C	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 [5]: display[display['UserId']=='AZY18LLTJ7J1NX']
```

Out[5]:

	UserId	ProductId	ProfileName	Time	Score	Text	COUNT(*)
80638	AZY10LLTJ7J1NX	B006P7E5ZJ	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 User-Id="AR5J8U146CURR"
ORDER BY ProductID
""", con)
display.head()
```

Out[7]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text
0	78445	B000HDL1RQ	AR5J8U146CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS	DELICIOUS WAFERS. I FIND THAT EUROPEAN WAFERS ...
1	138317	B000HDOPYC	AR5J8U146CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS	DELICIOUS WAFERS. I FIND THAT EUROPEAN WAFERS ...
2	138277	B000HDOPYM	AR5J8U146CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS	DELICIOUS WAFERS. I FIND THAT EUROPEAN WAFERS ...
3	73791	B000HDOPZG	AR5J8U146CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS	DELICIOUS WAFERS. I FIND THAT EUROPEAN WAFERS ...
4	155049	B000PAQ75C	AR5J8U146CURR	Geetha Krishnan	2	2	5	1199577600	LOACKER QUADRATINI VANILLA WAFERS	DELICIOUS WAFERS. I FIND THAT 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

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]: (62864, 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]: 89.80571428571429

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 calcaultions

```
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	Text
0	64422	B000MIDROQ	A161DK06JMCYF	J. E. Stephens "Jeanne"	3	1	5	1224992800	Bought This for My Son at College	My son loves spaghetti so I didn't hesitate or...
1	44737	B001EQG5RW	A2V09094FH7ABY	Ram	3	2	4	1212883200	Pure cocoa taste with crunchy almonds inside	It was almost a 'love at first bite' - the per...

```
In [12]: final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```



```
In [160]: tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=20, max_df=100)
tf_idf_vect.fit(X_tr)
print("some sample features(unique words in the corpus)",tf_idf_vect.get_feature_names()[0:10])
print("-"*50)

X_Tfidf_Tr = tf_idf_vect.transform(X_tr)
X_Tfidf_Cv = tf_idf_vect.transform(X_cv)
X_Tfidf_Test = tf_idf_vect.transform(X_test)

print("the type of count vectorizer ",type(X_Tfidf_Tr))
print("the shape of out text TFIDF vectorizer ",X_Tfidf_Tr.get_shape())
print("the number of unique words including both unigrams and bigrams ", X_Tfidf_Tr.get_shape()[1])

some sample features(unique words in the corpus) ['ability', 'able buy', 'able drink', 'able eat', 'able get', 'able make', 'able order', 'able purchase', 'able use', 'absolute best']
=====
the type of count vectorizer <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer (30882, 7134)
the number of unique words including both unigrams and bigrams 7134
```

[4.4] Word2Vec

```
In [40]: # Train your own Word2Vec model using your own text corpus
i=0
list_of_sentence=[]
for sentence in X_tr:
    list_of_sentence.append(sentence.split())

In [41]: # Using Google News Word2Vectors

# in this project we are using a pretrained model by google
# its 3.30 file, once you load this into your memory
# it occupies ~90b, so please do this step only if you have >12G of ram
# we will provide a pickle file with contains a dict ,
# and it contains all our corpus 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/007XKcUpISKDYNLNU7LSS2ipQwM/edit
# it's 1.90B in size.

# http://kavita-ganesan.com/gensim-word2vec-tutorial-starter-code/#.W175RFAz2PY
# 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=20,size=100, 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-negative300.bin', binary=True)
        #print(w2v_model.wv.most_similar('great'))
        #print(w2v_model.wv.most_similar('worst'))
    else:
        print("you don't have gopgle's word2vec file, keep want_to_train_w2v = True, to train your own w2v ")

=====

In [42]: w2v_words = list(w2v_model.wv.vocab)
print("number of words that occured minimum 5 times ",len(w2v_words))
print("sample words ", len(w2v_words))

number of words that occured minimum 5 times 4937
sample words 4937
```

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

```
In [43]: # average Word2Vec
# compute average word2vec for each review.

def getAvgWordToVector(list_of_sentence):
    sent_vectors = []; # the avg-w2v for each sentence/review is stored in this list
    for sentence in list_of_sentence: # for each review/sentence
        sent = sentence.split()
        sent_vec = np.zeros(100) # 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)
    return sent_vectors
```

[4.4.1.1] Avg W2v

```
In [44]: X_Avgw2V_Tr      = getAvgWordToVector(X_tr)

In [45]: X_Avgw2V_Cv      = getAvgWordToVector(X_cv)

In [46]: X_Avgw2V_Test    = getAvgWordToVector(X_test)
```

[4.4.1.2] TFIDF weighted W2v

```
In [47]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer(min_df=20, max_features=100)
tf_idf_matrix = model.fit(X_tr)
# 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_)))

In [53]: # 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

def getAvgw2vtfidfToVector(list_of_sentence):
    tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is stored in this list
    row=0;
    for sentence in list_of_sentence: # for each review/sentence
        sent = []
        sent_vec = np.zeros(100) # as word vectors are of zero length
        weight_sum =0; # num of words with a valid vector in the sentence/review
        sent = sentence.split()
        for word in sent: # for each word in a review/sentence3
            #print("word>>",word)
            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
    return tfidf_sent_vectors

In [54]: X_Avgw2vtfidf_Tr      = getAvgw2vtfidfToVector(X_tr)
X_Avgw2vtfidf_Cv      = getAvgw2vtfidfToVector(X_cv)
X_Avgw2vtfidf_Test      = getAvgw2vtfidfToVector(X_test)
```

FueatureEngineering :Adding length of each review, No of words as New Feature

In the Way of Improving the Accuracy we can do some feature Engineerings.

Here we are Adding up the Two New Feature to the Datasets.

- 1. No of Words in Each Review
- 2. The lenght of each Review

```
In [195]: def getNewFutre(old,org):
len_Tr      = []
wrd_Cnt_Tr = []
for i in org:
    l = []
    w = []
    l.append(len(i))
    len_Tr.append(l)
    w.append(len(i.split(' ')))
    wrd_Cnt_Tr.append(w)
old= np.hstack((old,len_Tr))
old= np.hstack((old,len_Tr))
return(old)

In [196]: X_Bow_Tr_new      = []
X_Bow_Cv_new      = []
X_Bow_Test_new      = []

X_Tfidf_Tr_new      = []
X_Tfidf_Test_new      = []
X_Tfidf_Cv_new      = []

X_Bow_Test_new      = getNewFutre(X_Bow_Test,X_test)
X_Bow_Tr_new      = getNewFutre(X_Bow_Tr,X_tr)
X_Bow_Cv_new      = getNewFutre(X_Bow_Cv,X_cv)

X_Tfidf_Tr_new      = getNewFutre(X_Tfidf_Tr,X_tr)
X_Tfidf_Test_new      = getNewFutre(X_Tfidf_Test,X_test)
X_Tfidf_Cv_new      = getNewFutre(X_Tfidf_Cv,X_cv)

In [197]: Bow_Feature.extend(["No of Words"," Lenght of Review"])

In [198]: tf_idf_Feature.extend(["No of Words"," Lenght of Review"])
```