

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

```
In [ ]: 1 # [1]. Reading Data
```

```
In [ ]: 1 ## [1.1] Loading the data
2
3 The dataset is available in two forms
4 1. .csv file
5 2. SQLite Database
6
7 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.
8 <br>
9
10 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 is positive and below 3, then the recommendation is negative.
```

```
In [1]: 1 %matplotlib inline
2 import warnings
3 warnings.filterwarnings("ignore")
4
5
6 import sqlite3
7 import pandas as pd
8 import numpy as np
9 import nltk
10 import string
11 import matplotlib.pyplot as plt
12 import seaborn as sns
13 from sklearn.feature_extraction.text import TfidfTransformer
14 from sklearn.feature_extraction.text import TfidfVectorizer
15
16 from sklearn.feature_extraction.text import CountVectorizer
17 from sklearn.metrics import confusion_matrix
18 from sklearn import metrics
19 from sklearn.metrics import roc_curve, auc
20 from nltk.stem.porter import PorterStemmer
21
22 import re
23 import string
24 from nltk.corpus import stopwords
25 from nltk.stem import PorterStemmer
26 from nltk.stem.wordnet import WordNetLemmatizer
27
28 from gensim.models import Word2Vec
29 from gensim.models import KeyedVectors
30 import pickle
31 import scipy as sp
32 from tqdm import tqdm
33 import os
34 from sklearn.cross_validation import train_test_split
35 from sklearn.neighbors import KNeighborsClassifier
36 from sklearn.metrics import accuracy_score
37 from sklearn.cross_validation import cross_val_score
38 from collections import Counter
39 from sklearn.metrics import accuracy_score
40 from sklearn import cross_validation
41 from sklearn.metrics import confusion_matrix
42 from sklearn.preprocessing import normalize
43 from sklearn import datasets, neighbors
44 from sklearn.metrics import roc_auc_score
45 from sklearn.preprocessing import StandardScaler
```

```
C:\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWarning: detected Windows; aliasing chunkize to chunkize_serial
warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")
```

```
C:\Anaconda3\lib\site-packages\sklearn\cross_validation.py:41: DeprecationWarning: This module was deprecated in version 0.18 in favor of the model_selection module into which all the refactored classes and functions are moved. Also note that the interface of the new CV iterators are different from that of this module. This module will be removed in 0.20.
"This module will be removed in 0.20.", DeprecationWarning)
```

```
In [2]: 1 # using SQLite Table to read data.
2 con = sqlite3.connect('C:\Python\Amazon review\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""", 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(1)
```

Number of data points in our data (525814, 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...

```
In [3]: 1 # [2] Exploratory Data Analysis
```

```
In [4]: 1 # 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:
2 #
3 # File "C:\Python\Amazon review\input-4-946a5d04d6db", line 3
4 # 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:
5 #
6 # SyntaxError: invalid syntax
```

```
In [5]: 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_position='last')
```

```
In [6]: 1 #Deduplication of entries
2 final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"}, keep='first', inplace=False)
3 final.shape
```

Out[6]: (364173, 10)

```
In [7]: 1 final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

[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 [8]: 1 # https://stackoverflow.com/questions/16206380/python-beautifulsoup-how-to-remove-all-tags-from-an-element
2 from bs4 import BeautifulSoup
3 # https://stackoverflow.com/a/47091490/4084039
4 import re
5
6 def decontracted(phrase):
7     # specific
8     phrase = re.sub(r"won't", "will not", phrase)
9     phrase = re.sub(r"can't", "can not", phrase)
10
11     # general
12     phrase = re.sub(r"n't", " not", phrase)
13     phrase = re.sub(r"re", " are", phrase)
14     phrase = re.sub(r"s", " is", phrase)
15     phrase = re.sub(r"d", " would", phrase)
16     phrase = re.sub(r"ll", " will", phrase)
17     phrase = re.sub(r"t", " not", phrase)
18     phrase = re.sub(r"ve", " have", phrase)
19     phrase = re.sub(r"m", " am", phrase)
20     return phrase
21
22 # https://gist.github.com/sebleier/554280
23 # we are removing the words from the stop words List: 'no', 'nor', 'not'
24 # <br /><br /> ==> after the above steps, we are getting "br br"
25 # we are including them into stop words List
26 # instead of <br /> if we have <br/> these tags would have remvomed in the 1st step
27
28 stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \
29 'you'll', "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
30 'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their', \
31 'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \
32 'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
33 'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'on', 'because', 'as', 'until', 'while', 'of', \
34 'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', \
35 'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', \
36 'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', \
37 'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
38 's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', \
39 've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', \
40 'hadn't', 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn', \
41 'mustn't', 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
42 'won', "won't", 'wouldn', "wouldn't"])
43
44 # Combining all the above stundents
45 from tqdm import tqdm
46 preprocessed_reviews = []
47 # tqdm is for printing the status bar
48 for sentence in tqdm(final['Text'].values):
49     sentence = re.sub(r"http\S+", "", sentence)
50     sentence = BeautifulSoup(sentence, 'lxml').get_text()
51     sentence = decontracted(sentence)
52     sentence = re.sub(r"\S*d\S+", "", sentence).strip()
53     sentence = re.sub(r"[^A-Za-z]+", ' ', sentence)
54     # https://gist.github.com/sebleier/554280
55     sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not in stopwords)
56     preprocessed_reviews.append(sentence.strip())
57
58 ## Similarly you can do preprocessing for review summary also.
59 preprocessed_sum = []
60 # tqdm is for printing the status bar
61 for sentence in tqdm(final['Summary'].values):
62     sentence = re.sub(r"http\S+", "", sentence)
63     sentence = BeautifulSoup(sentence, 'lxml').get_text()
64     sentence = decontracted(sentence)
65     sentence = re.sub(r"\S*d\S+", "", sentence).strip()
66     sentence = re.sub(r"[^A-Za-z]+", ' ', sentence)
67     # https://gist.github.com/sebleier/554280
68     sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not in stopwords)
69     preprocessed_sum.append(sentence.strip())
70
71 # adding the new preprocessed data as new columns to our final dataframe.
72
73 ps = pd.Series(preprocessed_sum)
74 final['Summary_new']=ps.values
75
76 pr = pd.Series(preprocessed_reviews)
77 final['Text_new']=pr.values
78
79 print('Shape of final',final.shape)
80 print(final['Score'].value_counts())
81 final.head(1)
82
83
84
85
86 # Saving the final data frame, prerocessed reviews and summary,
87 # so that we can resume directly from here without doing preprocessing again
88 # https://www.datacamp.com/community/tutorials/reading-writing-files-python
89
90 with open("final.txt", "wb") as file:
91     pickle.dump(final, file)
92
93 with open("preprocessed_reviews.txt", "wb") as file:
94     pickle.dump(preprocessed_reviews, file)
95
96 with open("preprocessed_summary.txt", "wb") as file:
97     pickle.dump(preprocessed_sum, file)
98
99
100 sorted_data = final.sort_values('Time', axis=0, ascending=True, inplace=False, kind='quicksort', na_position='last')
101
102 with open("final_sorted.txt", "wb") as file:
103     pickle.dump(sorted_data, file)
104
105 should probably open this file and pass the filehandle into BeautifulSoup.
106 BeautifulSoup." % markup)
107 92% [ ] | 335126/364171 [01:58<00:10, 2716.85it/s]C:\Anaconda3\lib\site-packages\bs4\_init_.py:219: UserWarning: "b'..." looks like a filename, not markup. You should probably open this file and pass the filehandle into BeautifulSoup.
108 BeautifulSoup." % markup)
109 97% [ ] | 354901/364171 [02:05<00:03, 2843.48it/s]C:\Anaconda3\lib\site-packages\bs4\_init_.py:219: UserWarning: "b'..." looks like a filename, not markup. You should probably open this file and pass the filehandle into BeautifulSoup.
110 BeautifulSoup." % markup)
111 99% [ ] | 358963/364171 [02:07<00:01, 2634.33it/s]C:\Anaconda3\lib\site-packages\bs4\_init_.py:219: UserWarning: "b'..." looks like a filename, not markup. You should probably open this file and pass the filehandle into BeautifulSoup.
112 BeautifulSoup." % markup)
113 99% [ ] | 359231/364171 [02:07<00:02, 2398.42it/s]C:\Anaconda3\lib\site-packages\bs4\_init_.py:219: UserWarning: "b'..." looks like a filename, not markup. You should probably open this file and pass the filehandle into BeautifulSoup.
114 BeautifulSoup." % markup)
115 100% [ ] | 364171/364171 [02:09<00:00, 2812.24it/s]

```

```
Shape of final (3641/1, 12)
1 307061
0 57110
Name: Score, dtype: int64
```

```
In [9]: 1 my_final = sorted_data[:100000]
2 my_final['Score'].value_counts()
```

```
Out[9]: 1 87729
0 12271
Name: Score, dtype: int64
```

```
In [10]: 1 from sklearn.model_selection import train_test_split
2 x = my_final['Text_new'].values
3 y = my_final['Score']
4
5 # Split the data set into train and test
6 X_1, X_test, y_1, Y_test = cross_validation.train_test_split(x, y, test_size=0.3, random_state=0)
7
8 # Split the train data set into cross validation train and cross validation test
9 X_train, X_cv, Y_train, Y_cv = cross_validation.train_test_split(X_1, y_1, test_size=0.3, random_state=0)
```

```
In [11]: 1 #https://stackoverflow.com/questions/10741346/numpy-most-efficient-frequency-counts-for-unique-values-in-an-array
2 # Printing the frequency of positive and negative values in Train, CV and Test data set
3 unique, counts = np.unique(Y_train, return_counts=True)
4
5 np.asarray((unique, counts)).T
```

```
Out[11]: array([[ 0, 6013],
[ 1, 42987]], dtype=int64)
```

```
In [12]: 1 unique, counts = np.unique(Y_test, return_counts=True)
2
3 np.asarray((unique, counts)).T
```

```
Out[12]: array([[ 0, 3665],
[ 1, 26335]], dtype=int64)
```

```
In [13]: 1 unique, counts = np.unique(Y_cv, return_counts=True)
2
3 np.asarray((unique, counts)).T
```

```
Out[13]: array([[ 0, 2593],
[ 1, 18407]], dtype=int64)
```

Bag of Words

```
In [14]: 1 # Please write all the code with proper documentation
2 # Bag of Words
3 count_vect = CountVectorizer() #in scikit-learn
4 train_bow = count_vect.fit_transform(X_train)
5 test_bow = count_vect.transform(X_test)
6 cv_bow = count_vect.transform(X_cv)
7
```

```
In [15]: 1 from scipy.sparse import save_npz
2 save_npz('train_bow.npz', train_bow)
```

```
In [16]: 1 save_npz('test_bow.npz', test_bow)
2 save_npz('cv_bow.npz', cv_bow)
```

Bag of Words with max_feature

```
In [26]: 1 # Please write all the code with proper documentation
2 # Bag of Words
3 count_vect = CountVectorizer(min_df=50, max_features=2000) #in scikit-learn
4 train_bowl = count_vect.fit_transform(X_train)
5 test_bowl = count_vect.transform(X_test)
6 cv_bowl = count_vect.transform(X_cv)
7
```

```
In [27]: 1 from scipy.sparse import save_npz
2 save_npz('train_bow_lim.npz', train_bowl)
```

```
In [28]: 1 save_npz('test_bow_lim.npz', test_bowl)
2 save_npz('cv_bow_lim.npz', cv_bowl)
```

TFIDF

```
In [17]: 1 # Please write all the code with proper documentation
2 tfidf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
3
4 final_tf_idf = tfidf_vect.fit_transform(X_train)
5
6 train_tfidf = tfidf_vect.fit_transform(X_train)
7 test_tfidf = tfidf_vect.transform(X_test)
8 cv_tfidf = tfidf_vect.transform(X_cv)
```

```
In [18]: 1 save_npz('train_tfidf.npz', train_tfidf)
2 save_npz('test_tfidf.npz', test_tfidf)
3 save_npz('cv_tfidf.npz', cv_tfidf)
```

TFIDF with max feature

```
In [29]: 1 # Please write all the code with proper documentation
2 tfidf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=50, max_features=2000)
3
4 final_tf_idf = tfidf_vect.fit_transform(X_train)
5
6 train_tfidf1 = tfidf_vect.fit_transform(X_train)
7 test_tfidf1 = tfidf_vect.transform(X_test)
8 cv_tfidf1 = tfidf_vect.transform(X_cv)
```

```
In [30]: 1 save_npz('train_tfidf_lim.npz', train_tfidf1)
2 save_npz('test_tfidf_lim.npz', test_tfidf1)
3 save_npz('cv_tfidf_lim.npz', cv_tfidf1)
```

W2V

```
# Train your own Word2Vec model using your own text corpus
1 i=0
2 list_of_sentence=[]
3 for sentence in X_train:
4     list_of_sentence.append(sentence.split())
5
6 # min_count = 5 considers only words that occurred atleast 5 times
7 w2v_model=Word2Vec(list_of_sentence,min_count=5,size=50, workers=4)
8
9 w2v_words = list(w2v_model.wv.vocab)
10 print("number of words that occurred minimum 5 times ",len(w2v_words))
11
12 def avgwtv(X_test):
13     '''
14     returns average word2vec
15     '''
16     i=0
17     list_of_sentence=[]
18     for sentence in X_test:
19         list_of_sentence.append(sentence.split())
20     test_vectors = []; # the avg-w2v for each sentence/review is stored in this list
21     for sent in tqdm(list_of_sentence): # for each review/sentence
22         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
23         cnt_words = 0; # num of words with a valid vector in the sentence/review
24         for word in sent: # for each word in a review/sentence
25             if word in w2v_words:
26                 vec = w2v_model.wv[word]
27                 sent_vec += vec
28                 cnt_words += 1
29             if cnt_words != 0:
30                 sent_vec /= cnt_words
31         test_vectors.append(sent_vec)
32     return test_vectors
33
34 train_avgw2v = avgwtv(X_train)
35 cv_avgw2v = avgwtv(X_cv)
36 test_avgw2v = avgwtv(X_test)
```

number of words that occurred minimum 5 times 13481

```
100% ██████████ 49000/49000 [02:02<00:00, 398.75it/s]
100% ██████████ 21000/21000 [00:55<00:00, 376.72it/s]
100% ██████████ 30000/30000 [01:19<00:00, 376.48it/s]
```

```
1 with open("train_avgw2v.txt", "wb") as file:
2     pickle.dump(train_avgw2v, file)
3 with open("cv_avgw2v.txt", "wb") as file:
4     pickle.dump(cv_avgw2v, file)
5 with open("test_avgw2v.txt", "wb") as file:
6     pickle.dump(test_avgw2v, file)
```

TFIDF W2V

```

1 # Please write all the code with proper documentation
2 # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
3 model = TfidfVectorizer()
4 model.fit(X_train)
5 # We are converting a dictionary with word as a key, and the idf as a value
6 dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
7
8 # TF-IDF weighted Word2Vec
9 tfidf_feat = model.get_feature_names() # tfidf words/col-names
10 # final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val = tfidf
11
12 #standardized_weight_w2v = StandardScaler().fit_transform(tfidf_sent_vectors)
13 #print(standardized_weight_w2v.shape)
14
15 def tfidfw2v(test):
16     ...
17     Returns tfidf word2vec
18     ...
19
20     tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is stored in this list
21     i=0
22     list_of_sentence=[]
23     for sentence in test:
24         list_of_sentence.append(sentence.split())
25
26     for sent in tqdm(list_of_sentence): # for each review/sentence
27         sent_vec = np.zeros(50) # as word vectors are of zero length
28         weight_sum = 0; # num of words with a valid vector in the sentence/review
29         for word in sent: # for each word in a review/sentence
30             if word in w2v_words and word in tfidf_feat:
31                 vec = w2v_model.wv[word]
32                 tf_idf = dictionary[word]*(sent.count(word)/len(sent))
33                 sent_vec += (vec * tf_idf)
34                 weight_sum += tf_idf
35             if weight_sum != 0:
36                 sent_vec /= weight_sum
37             tfidf_sent_vectors.append(sent_vec)
38
39     return tfidf_sent_vectors
40
41 train_tf2w = tfidfw2v(X_train)
42 cv_tf2w = tfidfw2v(X_cv)
43 test_tf2w = tfidfw2v(X_test)

```

```
1 with open("train_tf2v.txt", "wb") as file:
2     pickle.dump(train_tf2v, file)
3 with open("cv_tf2v.txt", "wb") as file:
4     pickle.dump(cv_tf2v, file)
5 with open("test_tf2v.txt", "wb") as file:
6     pickle.dump(test_tf2v, file)
```

```
In [24]: 1 with open("X_test.txt", "wb") as file:
2         pickle.dump(X_test, file)
3         with open("X_train.txt", "wb") as file:
4             pickle.dump(X_train, file)
5         with open("X_cv.txt", "wb") as file:
6             pickle.dump(X_cv, file)
```

```
In [25]: 1 with open("Y_test.txt", "wb") as file:
2         pickle.dump(Y_test, file)
3         with open("Y_train.txt", "wb") as file:
4             pickle.dump(Y_train, file)
5         with open("Y_cv.txt", "wb") as file:
6             pickle.dump(Y_cv, file)
```

```
In [ ]: 1
```