

Sentiment Analysis of Customer Reviews

Shivani Nandani

CS-306, Data Analysis and Visualization

Dhirubhai Ambani Institute of Information & Communication Technology

201801076@daiict.ac.in

1 Introduction

Sentiment analysis is the interpretation and classification of emotions (positive, negative and neutral) within text data using text analysis techniques. It helps gauge public opinion, conduct nuanced market research, monitor brand and product reputation, and understand customer experiences. Sentiment analysis, thus, has been an area of interest for the industry that relies heavily on consumer feedback. In this project, we have employed data analysis techniques on the Amazon Fine Food Reviews dataset. All the codes are uploaded on the the GitHub repository - [sentiment-analysis](#).

2 About the dataset

In this project, we have used the Amazon Fine Food Reviews dataset available on Kaggle. The dataset contains 568,454 reviews by 256,059 users for 74,258 products from Oct 1999 - Oct 2012. Around 260 users have more than 50 reviews.

Various parameter in the dataset are:

- Id - review number
- ProductId - unique identifier for the product
- UserId - unique identifier for the user
- ProfileName - profile name of the user
- HelpfulnessNumerator - number of users who found the review helpful
- HelpfulnessDenominator - number of users who indicated whether they found the review helpful or not
- Score - rating between 1 and 5
- Time - timestamp
- Summary - summary of the review
- Text - text of the review

3 Method and Observations

The method used in the project is as given below:

- Step 1 - Read data
- Step 2 - Data Preprocessing
 - check for missing values
 - check for duplicates
 - check for invalid entries
 - preprocessing text
 - * remove html tags
 - * remove alphanumeric
 - * convert to lowercase
 - * remove stopwords
 - * stemming the words
- Step 3 - Exploratory Data Analysis
 - distribution of score
 - review trend for each year
 - distribution of positive, negative and neutral reviews
 - understanding user data
 - understanding correlation

3.1 Read Data

To read data, we have used **Pandas** library. Output of the read is shown below:

Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text	
0	1	B00154F6B	A35DHTALM8G6	delmartian	1	1	5	1303862480	Good Quality Dog Food	I have bought several of the Vitality canned d...
1	2	B000135R6A	A10B7F4ZCVESXK	d11 pa	0	0	1	1146976000	Not as Advertised	Product arrived labeled as Jumbo Salted Peas...
2	3	B000LQ0CH0	ABRLM12XKA1N	Natalia Corres "Natalia Corres"	1	1	4	1119817600	"Delight" says it all	This is a confection that has been around a few...
3	4	B000UJ0Q1Q	A3958RC6FGUVV	Karl	3	3	2	1307923200	Cough Medicine	If you are looking for the secret ingredient i...
4	5	B000K22Z7K	A1UQ5C1F8G41T	Michael D. Bighae "M. Weisj"	0	0	5	1158777600	Great taffy at a great price. There was a wif...	
...
568449	568450	B001070H3D	A2805XK0D544Y	Lettie D. Carter	0	0	5	1299628800	Will not do without	Great for sesame chicken...this is a good if no...
568450	568451	B00051u1TCU	A13184VF0E81S	R. Sawyer	0	0	2	1131211200	disappointed	I'm disappointed with the flavor. The chocolate...
568451	568452	B0041613EE	A121A41QQ751Z	pked "pk_007"	2	2	5	1129782400	Perfect for our malipoos	These stars are small, so you can give 10-15 0...
568452	568453	B0041613EE	A318EVC7KXN0H	Kathy A. Welch "Lynell"	1	1	5	1131596800	Favorite Training and reward treat	These are the BEST treats for training and reu...
568453	568454	B000LRF2C02	A3LQ0PCZCVL9UC	srfe1117	0	0	5	1138422400	Great Honey	I am very satisfied ,product is as advertised,...
568454 rows x 10 columns										

568454 rows x 10 columns

Figure 1: df without data preprocessing

As is visible in the image, there are 568454 rows and 10 columns in the dataset.

3.2 Data Preprocessing

Data preprocessing is an essential part of data analysis. Since most real world problems are not bound in a specific boundary, datasets generally have incomplete data (missing values), noisy data, outliers or discrepancies. All these factors impact the success rate of a model, and it is hence important to perform data preprocessing to prepare it for further analysis.

3.2.1 Check for missing values

First, we check if any row has a missing value. In case it is present, we will drop that row. On running the test for missing value, we find that the dataset has no missing values.

3.2.2 Check for duplicates

Next, we check for duplicates in the dataset. Two data samples are considered same if they have the same values for UserID, ProfileName and Text.

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text
	3	4	8080640QJQ	A3580RCF8G0V	Karl	5	3	2 1307921200	Cough Medicine	If you are looking for the secret ingredient I...
	10	11	808081P8FE	A30K0C70W0R64	Canadian Fan	1	1	5 1307820800	The Best Hot Sauce in the World	I don't know if it's the catnip or the Tequila...
	20	30	8080918F8Y	A30K0C70W0R64	Canadian Fan	1	1	5 1307820800	The Best Hot Sauce in the World	I don't know if it's the catnip or the Tequila...
	60	09	808087VZ75	A312LAMB88AQ	calamence	0	0	3 1309251200	How much would you pay for a bag of chocolate like me, then \$0 is ok...	If you're sensitive like me, then \$0 is ok...
	69	70	808087VZ75	AWCF324C53F	C. Salcido	0	2	5 1307536000	pretzel haven!	this was sooooo delicious but too bad I ate it...

508409	508430	808030CLM4	A3P8BAG06PL7	Dark Kater Herold	3	3	3	1309051200	Quality & affordable Food	I was very pleased with the ingredient quality...
508410	508411	808030CLM4	A8BHLAC057MG	R2B	2	2	5	1332979200	Litter box	My main reason for the five star review was this...
508411	508412	808030CLM4	AUX3H58FX555	DAW	1	1	5	1315908000	Happy Camper	I bought this to try on two registered Maine Co...
508412	508413	808030CLM4	AU2305479Q0E8	Al Ling Chow	0	0	5	1336459200	Two Siberians like it!	when we brought home two 3-month-old bordered...
508413	508414	808030CLM4	A23Y20P4PMAL	k3m0s0e	1	2	2	1330041600	pristine edge cat food	My cats don't like it, what else can I say...
232828 rows x 10 columns										

Figure 2: df without duplicates found in df

We see that there are 232828 duplicated values. This includes all the versions of a duplicated data sample. Out of the multiple entries, we will retain the one with the maximum **HelpfulnessNumerator**. After appending the required samples, we get 393642 samples.

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text	
	0	1	8081646FS0	A3580RCF8G0V	delaertian	1	1	5	1303862400	Good Quality Dog Food	I have bought several of the vitality canned f...
	1	2	8081630R6A	A30B7F2CZVE9K	dill pa	0	1	1346970000	Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...	
	2	3	8080640QJQ	A8BHLAC057MG	Natalia Correa "Natalia Correa"	1	1	4	1219617000	"Delight" says it all	This is a confection that has been around a fe...
	4	5	808064227K	A30B8RCF8G0V	[Helen] D. Bigman "H. Wasiir"	0	5	1307927700	Great taffy at a great price. There are no...	Great taffy at a great price. There are no...	
	5	6	808064227K	A30B8RCF8G0V	Tuapanything	0	4	134261200	Nice Taffy	I got a wild hair for taffy and ordered this. I...	
	
508401	508402	808009YF01	A23RY20P4PMAL	Carlos Alvarez	0	0	1	1238081600	Terrible	These are just awful. I don't waste your money	
508402	508403	808030CLM4	A371P029129B	K. Bower	1	1	5	1275804000	Delicious!	I was pleasantly surprised when we bought this...	
508411	508412	808030CLM4	A30B8RCF8G0V	onthehookout	0	0	3	1343734000	TWINKERS CRACKERS DOG TREAT	My boxer loves this treat. He has not turned d...	
508412	508413	808030CLM4	A30B8RCF8G0V	Stone-Man	0	0	4	1337584000	Dog loves these!!	My dog goes crazy for these. They are deli...	
508406	508407	808030YF02	A30B8RCF8G0V	Lexia29	1	1	5	1344550000	Great product!!!	Fresh shipping, fresh unopened seeds, great per...	
393642 rows x 10 columns											

Figure 3: df without duplicates

3.2.3 Check for invalid entries

We will now check if the dataset has any invalid entries. An entry will be considered as invalid if the

HelpfulnessNumerator>**HelpfulnessDenominator**. We find two such entries as shown below:

Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text	
0	1	8081646FS0	A3580RCF8G0V	delaertian	1	1	5	1303862400	Good Quality Dog Food	I have bought several of the vitality canned f...
1	2	8081630R6A	A30B7F2CZVE9K	dill pa	0	0	1	1346970000	Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...
2	3	8080640QJQ	A8BHLAC057MG	Natalia Correa "Natalia Correa"	1	1	4	1219617000	"Delight" says it all	This is a confection that has been around a fe...
3	4	8080640QJQ	A30B8RCF8G0V	Karl	3	3	2	1307921200	Cough Medicine	If you are looking for the secret ingredient I...
4	5	808064227K	A30B8RCF8G0V	Michael D. Bigham "M. Bigham"	0	0	5	1307776000	Great taffy	Great taffy at a great price. There was a vid...
...
508400	508401	808030CLM4	A30B8RCF8G0V	Lettie D. Carter	0	0	5	1299628000	Will not do without	Great for sesame chicken...this is a good if no...
508400	508401	808030CLM4	A30B8RCF8G0V	Lettie D. Carter	0	0	2	1312511200	disappointed	I'm disappointed with the flavor...The chocolate...
508401	508402	808030CLM4	A23Y20P4PMAL	pkid "pkid"	2	2	5	1329782400	Perfect for our mutipoo	These stars are small, so you can give 10-15 o...
508402	508403	808030CLM4	A30B8RCF8G0V	Kathy A. Welch "kathel"	1	1	5	1315908000	Favorite Training and reward treat	These are the BEST treats for training and reward...
508403	508404	808030CLM4	A30B8RCF8G0V	srfe1117	0	0	5	1338422400	Great Honey	I am very satisfied...product is as advertised...
508404 rows x 10 columns										

Figure 4: df without duplicates

On removing these two entries, we get 393640 entries, which is then re-indexed to match the new state.

#	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text
0	1	8081646FS0	A3580RCF8G0V	delaertian	1	1	5	1303862400	Good Quality Dog Food	I have bought several of the vitality canned f... d...
1	2	8081630R6A	A30B7F2CZVE9K	dill pa	0	0	1	1346970000	Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...
2	3	8080640QJQ	A8BHLAC057MG	Natalia Correa "Natalia Correa"	1	1	4	1219617000	"Delight" says it all	This is a confection that has been around a fe...
3	4	8080640QJQ	A30B8RCF8G0V	Karl	3	3	2	1307921200	Cough Medicine	If you are looking for the secret ingredient I...
4	5	808064227K	A30B8RCF8G0V	Michael D. Bigham "M. Bigham"	0	0	5	1307776000	Great taffy	Great taffy at a great price. There was a vid...
...
393639	508400	808030CLM4	A30B8RCF8G0V	Lettie D. Carter	0	0	5	1299628000	Will not do without	Great for sesame chicken...this is a good if no...
393640	508401	808030CLM4	A30B8RCF8G0V	Lettie D. Carter	0	0	2	1312511200	disappointed	I'm disappointed with the flavor...The chocolate...
393641	508402	808030CLM4	A23Y20P4PMAL	pkid "pkid"	2	2	5	1329782400	Perfect for our mutipoo	These stars are small, so you can give 10-15 o...
393642	508403	808030CLM4	A30B8RCF8G0V	Kathy A. Welch "kathel"	1	1	5	1315908000	Favorite Training and reward treat	These are the BEST treats for training and reward...
393643	508404	808030CLM4	A30B8RCF8G0V	srfe1117	0	0	5	1338422400	Great Honey	I am very satisfied...product is as advertised...
393640 rows x 10 columns										

Figure 5: df without duplicates and invalid entries

3.2.4 Preprocessing text

For preprocessing the text (i.e., the **Review**), we do the following:

- remove html tags
- remove alphanumeric
- convert to lowercase
- remove stopwords
- stemming the words

HTML tags are introduced when the user adds links other such components to the review. Alphanumeric characters (such as numbers, special characters etc.) are also removed as they do not add any *sentimental value* to the review. To correctly judge the words we change all the letters to lowercase so that they can all be treated equally, without creating unnecessary duplicates. Finally, we deal with stopwords and stemming of words. Stopwords are the words that do not add any meaning to the sentence and thus can be ignored without sacrificing the essence of the sentence. Examples of stopwords are words like *the*, *are*, *have* etc. Stemming of words involves reducing the words to their base form. For example, on stemming,

The final dataset after complete data preprocessing has 393640 values as shown below:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text
0	1	B001E4CF00	AJ55C8W40R0G	delmarjian	1	1	5	1303862400	Good Quality Dog Food	bought sever vital dog food product found good.

```
393640 rows x 10 columns
```

Figure 6: df after complete data preprocessing

3.3 Exploratory Data Analysis

3.3.1 Distribution of scores

From the below-given histogram, we can see that most reviews are positive.

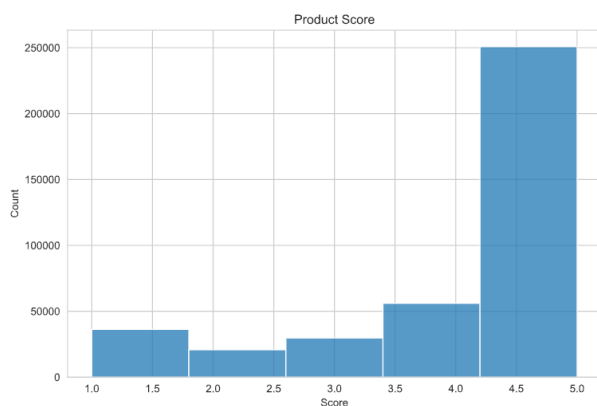


Figure 7: score distribution among various reviews

3.3.2 Review trend for each year

Here we can see that in the initial years, i.e., from 2001 to 2006, the number of reviews remain almost constant. However, the rate increases after 2006 and most reviews after that period are positive.

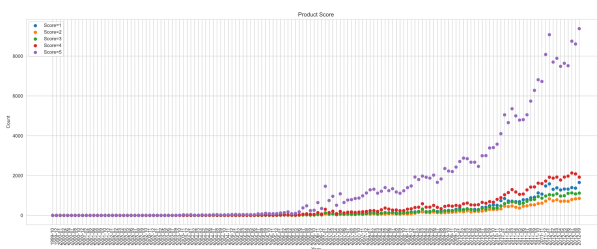


Figure 8: review trend for each year (month-wise manner)

3.3.3 Distribution of positive, negative and neutral reviews

We divided the reviews in three categories based on their score – any review with score greater than 3 has be considered positive; any review with score equal to 3 is consider neutral; any review with score less than 3 has be considered negative.

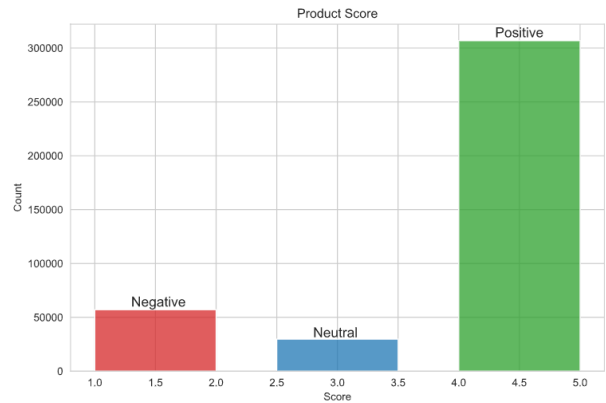


Figure 9: histogram showing number of positive, negative and neutral reviews

3.3.4 Understanding user data

We can see that there are 256055 unique users. Since the mean is 1.537326, we can say that most users bought the product only once.

```
count      256055.000000
mean         1.537326
std          2.732791
min          1.000000
25%          1.000000
50%          1.000000
75%          1.000000
max          329.000000
Name: UserId, dtype: float64
```

Figure 10: user data

3.3.5 Wordclouds

Wordclouds are display the most frequently occurring words in the given sample. Here, the size of a word is proportional to its frequency.



Figure 11: wordcloud for all reviews



Figure 12: wordcloud for positive reviews



Figure 13: wordcloud for neutral reviews



Figure 14: wordcloud for negative reviews

3.3.6 Understanding correlation

We checked the correlation between various fields:

Parameter 1	Parameter 2	Correlation
Score	ProductID	-0.03879
Score	UserId	-0.0139827
Score	HelpfulnessNumerator	-0.03620
Score	Length of the review	-0.05968

4 Conclusions

- In FIG[7] and FIG[9], we see that most reviews are positive. This can be due to the reasons discussed below. However, a consequence of this is that the dataset is not balanced. Thus, in this case, if a model to predict if a review is positive or negative is trained using this dataset, it will be inherently biased towards positive reviews. In

such cases, accuracy will give a false sense of correctness of the model.

- In FIG[8], we see that most reviews are positive after 2006. This can be interpreted in two ways:

1. The number of reviews grew with the increases number of users of Amazon, the graph of which is give below¹. The nature of this graph is similar to that of FIG[8].

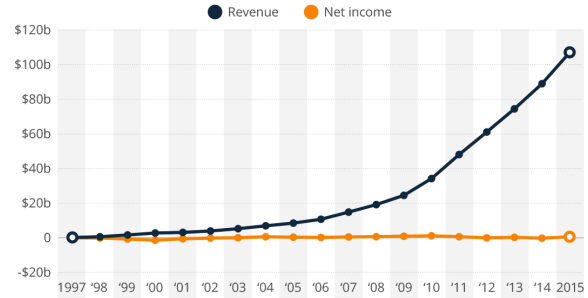


Figure 15: growth of Amazon

2. This trend seen in FIG[8] can also be explained by the marketing tactics used by the companies in which the users are either paid to write fake positive reviews, or company employees write positive reviews using fake accounts.

However, the actual reason can be a combination of above-mentioned reasons.

- In FIG[10], we see that the maximum number of reviews posted by a single user is 329. This *UserId* acts as an outlier in this case since it is extremely far from the mean and variance of the data. In real life context, this can be related to the one of the reasons for the sudden growth mentioned above. We checked the distribution of scores for this users, we find that most of the reviews are rated 4. This, thus, can be a proof of the fraudulent ways in which companies try to increase their customer ratings.

	Id	ProductId	HelpfulnessNumerator	HelpfulnessDenominator	Score
count	329.000000	329.000000	329.000000	329.000000	329.000000
mean	270291.358663	38124.367781	1.124620	1.234043	3.659574
std	170106.232415	19451.560950	2.728514	2.968831	0.657676
min	110.000000	17.000000	0.000000	0.000000	1.000000
25%	99740.000000	22090.000000	0.000000	0.000000	3.000000
50%	271766.000000	40085.000000	0.000000	0.000000	4.000000
75%	409657.000000	54438.000000	1.000000	1.000000	4.000000
max	566680.000000	67248.000000	20.000000	22.000000	4.000000

Figure 16: details for UserId with 329 reviews

- From FIG[11], FIG[12], FIG[??] and FIG[14], we see that most words associated with the complete set are similar to the ones associated with positive and neutral sets. This can be attributed

¹Source - [One simple chart that shows Amazon's relentless focus on long-term growth](#)

to the imbalanced dataset. In the wordcloud for negative set, we see words such as good and great. This is because they were used with the word 'not' to imply dissatisfaction. However, these were removed as a part of text preprocessing during the removal of stopwords. Finally, we see that the negative set contains a lot of conditional words (such as 'however', 'though' etc.). This can show that most users were not completely dissatisfied with the product but the same did not match their expectations.

- From the correlation data we can see that the parameters are not dependent on each other. Thus, an individual doesn't always rate a product good or bad, but do so based on their experiences with the products. Similarly, most products have their rating fluctuating. Even with *Sentiment* as a parameter instead of *Score*², we see that correlation does not improve.

References

- [1] [Sentiment Analysis: Concept, Analysis and Applications](#)
- [2] [Sentiment Analysis Explained](#)
- [3] [Techniques and applications for sentiment analysis](#)
- [4] [A Gentle Introduction to Imbalanced Classification](#)
- [5] [Learning from imbalanced data.](#)

²given in the GitHub repository