

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.

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	B001E4K6GB	A352M7ALP86W	delmartin	1	1	5	1303862400	Good Quality Dog Food	I have bought several of the Vitality canned d...
1	2	B00B13JRG4	A1D8FFZCZV5NK	d11 pa	0	0	1	1346976000	Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...
2	3	B00OLQ0CH8	ABRLN2J2X41N	Natalia Correa "Natalia Correa"	1	1	4	1219827600	"Delight" says it all	This is a confection that has been around a fe...
3	4	B000UAHQDQ	A395B0RC5FQ0V	Karl	3	3	2	1307923200	Cough Medicine	If you are looking for the secret ingredient I...
4	5	B006K2Z27K	A1UQ5CLF8G4T	Michael D. Biglow "M. Weis"	0	0	5	1350777600	Great taffy	Great taffy at a great price. There was a wid...
...
568449	568450	B001E07N8B	A2B85XOR0544Y	Lettie D. Carter	0	0	5	1299628800	Will not do without	Great for sesame chicken...this is a good if no...
568450	568451	B0005UWTCU	A13BAFVPE8K15	R. Sawyer	0	0	2	1312512000	disappointed	I'm disappointed with the flavor. The chocolate...
568451	568452	B0041613EE	A123AASQV751Z	pked "pk_007"	2	2	5	1329782400	Perfect for our malisipo	These stars are small, so you can give 10-15 o...
568452	568453	B0041613EE	A13BEVCTKXNDH	Kathy A. Welch "Kamel"	1	1	5	1315596800	Favorites Training and reward treat	These are the BEST treats for training and rew...
568453	568454	B001LR2C02	A3LQ0PJCZVL9C	srfa1117	0	0	5	1338422400	Great Honey	I am very satisfied product is as advertised...

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
5	4	B00N8ANQDQ	A39B0R0CF0XV	Karl	3	3	2	130792300	Cough Medicine
10	11	B00N8PFF9F	A3H0K70AMQK4	Canadian Fan	1	1	5	1307820800	The Best Hot Sauce in the world
29	30	B00N8PFF9F	A3H0K70AMQK4	Canadian Fan	1	1	5	1307820800	The Best Hot Sauce in the world
68	69	B00N8P7V7S	A3KL2L6AB880Q2	calmesense	0	0	5	1309251200	How much would you pay for a bag of chocolate
69	70	B00N8P7V7S	A3KBF2Z6367F	C. Salido	0	2	5	1305753600	pretzel haven!
...
508409	508410	B00N8LCU8A	A3F88AG0UPL7	Dark Water Herald	3	3	5	1306051200	Quality & affordable food
508410	508411	B00N8LCU8A	A3B8LAC057M6	R20	2	2	5	1332979200	Litter box
508411	508412	B00N8LCU8A	A3U8H78FK55S	DAN	1	1	5	1319508800	Happy Camper
508412	508413	B00N8LCU8A	A3V20D47PQ8E	Al Chow	0	0	5	1336452800	Two Siberians like it!
508413	508414	B00N8LCU8A	A3ZY28K7P8AL	klascobe	1	2	2	1330041600	premium edge cat food

232828 rows x 10 columns

Figure 2: all duplicates found in df

We see that there are 232828 duplicates values. This includes all the versions of a duplicated data sample. Out of the multiple entries, we will retain the with 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	B00E248F5Q	A35GH7AH0R0W	debertart	1	1	5	1303862400	Good Quality Dog Food
1	2	B00B830G6A	A3D87FZC2V59K	dll pa	0	0	1	1346976000	Not as Advertised
2	3	B00N8Q0CH0	A3BLN61Z0X03N	Natalia Corres "Natalia Corres"	1	1	4	1219617600	"Delight" says it all
4	5	B00K6Z2Z7K	A3U9C5L7P8QJ7	Michael D. Bigham "M. Basmir"	0	0	5	1350777600	Great taffy
5	6	B00K6Z2Z7K	A39B0R0CF0XV	Tuesapenytling	0	0	4	1342051200	Nice Taffy
...
508041	508042	B00N8P7E2Y	A3R9K0ZC8B5H	Carlos Alvarez	0	0	1	1239081600	Terrific, don't waste your money
508042	508043	B00N82059N	A3T7P0ZK390B	K. Bauer	1	1	5	1275004800	Delicious!
509111	509112	B00N87L8UW	A3K28ED7P8M5	anthelabout	0	0	3	1343174400	THESEERS CHICKEN SOU BEST
509112	509113	B00N87L8UW	A3AF8B0C358UQ2	Stone-Man	0	0	4	1337558400	My dog goes crazy for these. They are BIG like...
500004	500007	B00N85VQ18	A3U8H78FK55S	Lexie29	1	1	5	1344556800	Great product!!!

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	B00E248F5Q	A35GH7AH0R0W	debertart	1	1	5	1303862400	I have bought several of the Victory brand dog food...
1	2	B00B830G6A	A3D87FZC2V59K	dll pa	0	0	1	1346976000	Not as advertised
2	3	B00N8Q0CH0	A3BLN61Z0X03N	Natalia Corres "Natalia Corres"	1	1	4	1219617600	"Delight" says it all
3	4	B00N8ANQDQ	A39B0R0CF0XV	Karl	3	3	2	1307923000	Cough Medicine
4	5	B00K6Z2Z7K	A3U9C5L7P8QJ7	Michael D. Bigham "M. Basmir"	0	0	5	1350777600	Great taffy at a great price. There was a mild...
...
508409	508410	B00N8LCU8A	A3F88AG0UPL7	Lettie D. Carter	0	0	5	1299628800	Will not do without
508410	508411	B00N8LCU8A	A3B8LAC057M6	R. Sawyer	0	0	2	1331251200	disappointed
508411	508412	B00N8LC1EE	A321A41Q0751Z	pkid "pk_m07"	2	2	5	1320782400	Perfect for our medium sized dogs
508412	508413	B00N8LC1EE	A32REV7K0N0H	Kathy A. Welch "kathel"	1	1	5	1331596800	Favorite Training and reward treat
508413	508414	B00N8LC20E	A3U9C5L7P8QJ7	vr4all17	0	0	5	1336232000	Great Honey

508454 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	B00E248F5Q	A35GH7AH0R0W	debertart	1	1	5	1303862400	Good Quality Dog Food	
1	2	B00B830G6A	A3D87FZC2V59K	dll pa	0	0	1	1346976000	Product arrived labeled as Jumbo Salted Peanut...	
2	3	B00N8Q0CH0	A3BLN61Z0X03N	Natalia Corres "Natalia Corres"	1	1	4	1219617600	"Delight" says it all	
3	4	B00N8ANQDQ	A39B0R0CF0XV	Karl	3	3	2	130792300	Cough Medicine	
4	5	B00K6Z2Z7K	A3U9C5L7P8QJ7	Michael D. Bigham "M. Basmir"	0	0	5	1350777600	Great taffy at a great price. There was a vid...	
...	
393635	508410	B00N8LCU8A	A3F88AG0UPL7	Dark Water Herald	3	0	0	5	1299628800	Will not do without
393636	508411	B00N8LCU8A	A3B8LAC057M6	R. Sawyer	0	0	2	1331251200	disappointed	
393637	508412	B00N8LC1EE	A321A41Q0751Z	pkid "pk_m07"	2	2	5	1320782400	Perfect for our malinois	
393638	508413	B00N8LC1EE	A32REV7K0N0H	Kathy A. Welch "kathel"	1	1	5	1331596800	Favorite treats for training and reward	
393639	508414	B00N8LC20E	A3U9C5L7P8QJ7	vr4all17	0	0	5	1336232000	Great Money	
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, *eating* and *eats* are both changed to *eat*. Stemming is an important part of text preprocessing as it gives

us the root for each value, which allows the model to consider all valid versions of the word.

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	B001EAF6P	A35DKTAMURGH	delmarian	1	1	5	1303862400	Good Quality Dog Food bought sever vital dog food product found good...
1	2	B00815G54	A3D8T4C2E9SE	dill pa	0	0	1	1346970800	Not as Advertised product arriv label jumbo salt peanut peanut a...
2	3	B000LQ0CH	ABRLNCTKXAZN	Rutalia Correa "Rutalia Correa"	1	1	4	1219617600	"Delight" says It all confect around centuri lighte pinked citrea ge...
3	4	B000UAPQD	A395BRCF60NV	Karl	3	3	2	1307923200	Cough Medicine Jock sacre ingredi robitussin halter food ge...
4	5	B00K6Z2Z7E	A3UQ5CLF8G4T	Michael D. Bigham "Bigham M. Bigham"	0	0	5	1358777600	Great taffy great taffi great price ude assort yams taff...
...
393635	568450	B001E0702D	A28WZKOR564Y	Lettie D. Carter	0	0	5	1259628800	Will not do without great sasem chicken good better reclar waten h...
393636	568451	B00353ATCU	A35BAFVFE8K25	R. Sawyer	0	0	2	1311291200	disappointed disappoint flavor choco some aspect weak milk...
393637	568452	B00451LIEE	A21AASQ07512	pked "pk_n07"	2	2	5	1320784800	Perfect for our puppies star small give 10 15 new train session tri ty...
393638	568453	B00451LIEE	A32BEVTKX0NH	Kathy A. Walsh "Kathy A. Walsh"	1	1	5	1311596800	Favorite Training and reward treat best treat train reward dig good groom lower s...
393639	568454	B001LKZCJ2	A3UQ5CLF8G4T	srffell17	0	0	5	1338422400	Great Honey satisfi product adverti use cereal new vinge...

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.

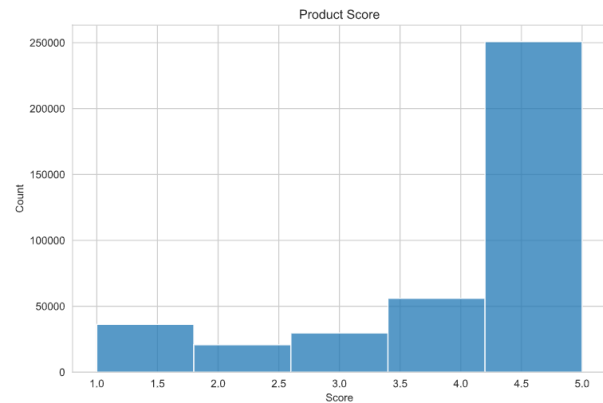


Figure 7: score distribution among various reviews

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.

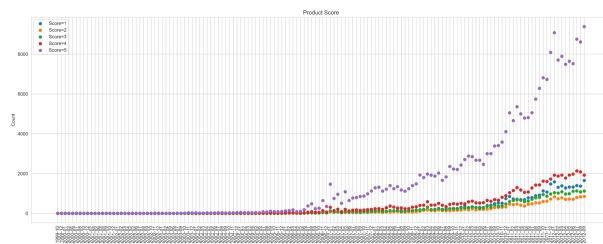


Figure 8: review trend for each year (month-wise manner)

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

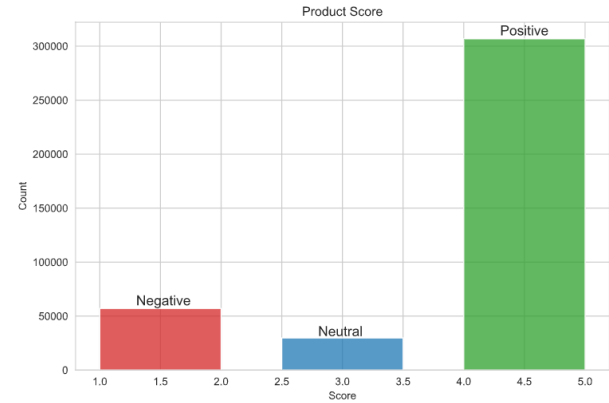


Figure 9: histogram showing number of positive, negative and neutral reviews

4 Conclusions

- write about how the dataset is unbalanced and this will hinder the modeling process
- In FIG[8], we see that most reviews are positive after 2006. This can be interpreted in two ways:

- 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].

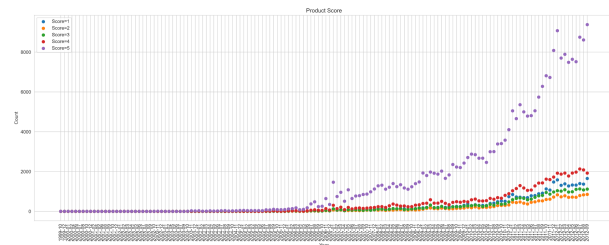


Figure 10: review trend for each year (month-wise manner)

- This trend seen is FiG[8] can also be explained by the marketing tactics used by

¹Source - One simple chart that shows Amazon's relentless focus on long-term growth

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.

References

Please number citations consecutively within brackets [1]. The sentence punctuation follows the bracket [2]. Refer simply to the reference number, as in [3]—do not use “Ref. [3]” or “reference [3]” except at the beginning of a sentence: “Reference [3] was the first ...”

Number footnotes separately in superscripts. Place the actual footnote at the bottom of the column in which it was cited. Do not put footnotes in the abstract or reference list. Use letters for table footnotes.

Unless there are six authors or more give all authors’ names; do not use “et al.”. Papers that have not been published, even if they have been submitted for publication, should be cited as “unpublished” [4]. Papers that have been accepted for publication should be cited as “in press” [5]. Capitalize only the first word in a paper title, except for proper nouns and element symbols.

For papers published in translation journals, please give the English citation first, followed by the original foreign-language citation.

References

- [1] G. Eason, B. Noble, and I. N. Sneddon, “On certain integrals of Lipschitz-Hankel type involving products of Bessel functions,” *Phil. Trans. Roy. Soc. London*, vol. A247, pp. 529–551, April 1955.
- [2] J. Clerk Maxwell, *A Treatise on Electricity and Magnetism*, 3rd ed., vol. 2. Oxford: Clarendon, 1892, pp.68–73.
- [3] I. S. Jacobs and C. P. Bean, “Fine particles, thin films and exchange anisotropy,” in *Magnetism*, vol. III, G. T. Rado and H. Suhl, Eds. New York: Academic, 1963, pp. 271–350.
- [4] K. Elissa, “Title of paper if known,” unpublished.
- [5] R. Nicole, “Title of paper with only first word capitalized,” *J. Name Stand. Abbrev.*, in press.
- [6] M. Young, *The Technical Writer’s Handbook*. Mill Valley, CA: University Science, 1989.