



Ratings Project

Submitted by:
ASHEEM SIWACH

ACKNOWLEDGMENT

I would like to thank Flip Robo Technologies for providing me with the opportunity to work on this project from which I have learned a lot. I am also grateful to Mr. Shubham Yadav for his constant guidance and support.

Reference sources are :-

- Google
- Stackoverflow.com
- Analytics vidhya
- Notes and repository from DataTrained

INTRODUCTION

• Business Problem Framing

We have a client who has a website where people write different reviews for technical products. Now they are adding a new feature to their website i.e. The reviewer will have to add stars(rating) as well with the review. The rating is out 5 stars and it only has 5 options available 1 star, 2 stars, 3 stars, 4 stars, 5 stars. Now they want to predict ratings for the reviews which were written in the past and they don't have a rating. So, we have to build an application which can predict the rating by seeing the review.

• Conceptual Background of the Domain Problem

The advent of electronic commerce with growth in internet and network technologies has led customers to move to online retail platforms such as Amazon, Walmart, Flip Kart, etc. People often rely on customer reviews of products before they buy online. These reviews are often rich in information describing the product. Customers often choose to compare between various products and brands based on whether an item has a positive or negative review. More often, these reviews act as a feedback mechanism for

the seller. Through this medium, sellers strategize their future sales and product improvement.

- Motivation for the Problem Undertaken

There are many websites/applications who are providing online shopping services. But as there is no direct contact between the seller and consumer it becomes somewhat difficult for the seller what customers want and what they should provide to their customers.

Every day we come across various products in our lives, on the digital medium we swipe across hundreds of product choices under one category. It will be tedious for the customer to make selection. Here comes 'reviews' where customers who have already got that product leave a rating after using them and brief their experience by giving reviews.

As we know ratings can be easily sorted and judged whether a product is good or bad. But when it comes to sentence reviews, we need to read through every line to make sure the review conveys a positive or negative sense. In the era of artificial intelligence, things like that have got easy with the Natural Language Processing (NLP) technique. Therefore, it is important to minimize the number of false positives our model produces, to encourage all constructive conversation.

Our model also provides beneficence for the platform hosts as it replaces the need to manually moderate discussions, saving time and resources. Employing a machine learning model to predict ratings promotes easier way to distinguish between products qualities, costs and many other features.

Analytical Problem Framing

- Data Collection

For data collection we have used Selenium Web Scrapping tool. By using this tool we gathered almost 20000 rows data which included of smartphones, watches, monitors, refrigerators, headphones, microwaves etc.

Together we collected data and then combine all into a single dataframe using pandas dataframe.

```
[2]: 1 # Calling webdriver and going to the url
2 driver = webdriver.Chrome("chromedriver.exe")
3
4 driver.get('https://www.amazon.in/')

[3]: 1 # Using search bar for finding the desired data to scrape for smartphones
2 driver.find_element_by_xpath("//input[@id='twotabsearchtestbox']").send_keys('smartphones',Keys.ENTER)
3

[4]: 1 # Create empty list to store product urls and then append url into it
2 product_url = []
3 for i in tqdm_notebook(range(0,2)):
4     URLLS = driver.find_elements_by_xpath("//div[@class='a-section a-spacing-none']/div/h2/a")
5     for u in URLLS:
6         product_url.append(u.get_attribute("href"))
7
8     # Going to next page for more url
9     try:driver.find_element_by_xpath("//*[@id='search']/div[1]/div/div[1]/div/span[3]/div[2]/div[20]/span/div/div/ul/li[7]/a")
10    except NoSuchElementException:
11        pass
```

```
In [10]: 1 # Using forloop to scrape the data from each product url
2 for url in tqdm_notebook(product_url):
3     driver.get(url)
4     time.sleep(1.5)
5
6     # going to ratings page
7     try:driver.find_element_by_xpath("//*[@id='acrCustomerReviewLink']").click()
8     except NoSuchElementException:
9         print(cl("No ratings available",color='blue',attrs=['dark']))
10        pass
11
12    # Click to see all reviews of product
13    try:driver.find_element_by_xpath("//*[@id='reviews-medley-footer']/div[2]/a").click()
14    except NoSuchElementException:
15        print(cl("No more reviews",color='green',attrs=['dark']))
16        pass
17
18    time.sleep(1.5)
19    for i in tqdm_notebook(range(0,30)):
20        time.sleep(1)
21
22        # fetch ratings of product from each review
23        try:
24            RATINGS = driver.find_elements_by_xpath("//div[@class='a-section celwidget']/div[2]/a[1]")
25            for r in RATINGS:
26                driver.implicitly_wait(2)
27
28                # fetching reviews of product
29                try:
30                    REVIEWS = driver.find_elements_by_xpath("//div[@class='a-section celwidget']/div[4]")
31                    for v in REVIEWS:
32                        reviews.append(v.text.replace('\n',''))
33                except NoSuchElementException:
34                    reviews.append("Not Available")
35
36                #Printing the page scrapped
37                print("Scraping page {} ".format(i+1))
38
39                # Going to next page of reviews
40                try:driver.find_element_by_xpath("//*[@id='cm_cr-pagination_bar']/ul/li[2]/a").click()
41                except NoSuchElementException:
42                    print("No more reviews")
43                    pass
```

```

In [13]: 1 # Storing smartphones scraped data into dataframe
          2 smartphones=pd.DataFrame({"Reviews":reviews,"Ratings":ratings})
          3 smartphones.to_csv("smartphones.csv")

In [14]: 1 # Storing data into product dataframe
          2 products=pd.DataFrame({"Reviews":reviews,"Ratings":ratings})
          3 products.to_csv("products.csv")

In [15]: 1 # Storing scraped data of laptops into product dataframe
          2 laptops=pd.DataFrame({"Reviews":reviews,"Ratings":ratings})
          3 # Storing data into csv file
          4 laptops.to_csv("laptops.csv")

In [17]: 1 # Storing reviews and rating data of headphones
          2 headphones=pd.DataFrame({"Reviews":reviews[:4000],"Ratings":ratings[:4000]})
          3 headphones

Out[17]:
Reviews  Ratings
0         1         4
1         1         4
2         1         4
3         1         4
4         1         4
5         1         4
6         1         4
7         1         4
8         1         4
9         1         4
10        1         4
11        1         4
12        1         4
13        1         4
14        1         4
15        1         4
16        1         4
17        1         4
18        1         4
19        1         4
20        1         4
21        1         4
22        1         4
23        1         4
24        1         4
25        1         4
26        1         4
27        1         4
28        1         4
29        1         4
30        1         4
31        1         4
32        1         4
33        1         4
34        1         4
35        1         4
36        1         4
37        1         4
38        1         4
39        1         4
40        1         4
41        1         4
42        1         4
43        1         4
44        1         4
45        1         4
46        1         4
47        1         4
48        1         4
49        1         4
50        1         4
51        1         4
52        1         4
53        1         4
54        1         4
55        1         4
56        1         4
57        1         4
58        1         4
59        1         4
60        1         4
61        1         4
62        1         4
63        1         4
64        1         4
65        1         4
66        1         4
67        1         4
68        1         4
69        1         4
70        1         4
71        1         4
72        1         4
73        1         4
74        1         4
75        1         4
76        1         4
77        1         4
78        1         4
79        1         4
80        1         4
81        1         4
82        1         4
83        1         4
84        1         4
85        1         4
86        1         4
87        1         4
88        1         4
89        1         4
90        1         4
91        1         4
92        1         4
93        1         4
94        1         4
95        1         4
96        1         4
97        1         4
98        1         4
99        1         4
100       1         4
101       1         4
102       1         4
103       1         4
104       1         4
105       1         4
106       1         4
107       1         4
108       1         4
109       1         4
110       1         4
111       1         4
112       1         4
113       1         4
114       1         4
115       1         4
116       1         4
117       1         4
118       1         4
119       1         4
120       1         4
121       1         4
122       1         4
123       1         4
124       1         4
125       1         4
126       1         4
127       1         4
128       1         4
129       1         4
130       1         4
131       1         4
132       1         4
133       1         4
134       1         4
135       1         4
136       1         4
137       1         4
138       1         4
139       1         4
140       1         4
141       1         4
142       1         4
143       1         4
144       1         4
145       1         4
146       1         4
147       1         4
148       1         4
149       1         4
150       1         4
151       1         4
152       1         4
153       1         4
154       1         4
155       1         4
156       1         4
157       1         4
158       1         4
159       1         4
160       1         4
161       1         4
162       1         4
163       1         4
164       1         4
165       1         4
166       1         4
167       1         4
168       1         4
169       1         4
170       1         4
171       1         4
172       1         4
173       1         4
174       1         4
175       1         4
176       1         4
177       1         4
178       1         4
179       1         4
180       1         4
181       1         4
182       1         4
183       1         4
184       1         4
185       1         4
186       1         4
187       1         4
188       1         4
189       1         4
190       1         4
191       1         4
192       1         4
193       1         4
194       1         4
195       1         4
196       1         4
197       1         4
198       1         4
199       1         4
200       1         4
201       1         4
202       1         4
203       1         4
204       1         4
205       1         4
206       1         4
207       1         4
208       1         4
209       1         4
210       1         4
211       1         4
212       1         4
213       1         4
214       1         4
215       1         4
216       1         4
217       1         4
218       1         4
219       1         4
220       1         4
221       1         4
222       1         4
223       1         4
224       1         4
225       1         4
226       1         4
227       1         4
228       1         4
229       1         4
230       1         4
231       1         4
232       1         4
233       1         4
234       1         4
235       1         4
236       1         4
237       1         4
238       1         4
239       1         4
240       1         4
241       1         4
242       1         4
243       1         4
244       1         4
245       1         4
246       1         4
247       1         4
248       1         4
249       1         4
250       1         4
251       1         4
252       1         4
253       1         4
254       1         4
255       1         4
256       1         4
257       1         4
258       1         4
259       1         4
260       1         4
261       1         4
262       1         4
263       1         4
264       1         4
265       1         4
266       1         4
267       1         4
268       1         4
269       1         4
270       1         4
271       1         4
272       1         4
273       1         4
274       1         4
275       1         4
276       1         4
277       1         4
278       1         4
279       1         4
280       1         4
281       1         4
282       1         4
283       1         4
284       1         4
285       1         4
286       1         4
287       1         4
288       1         4
289       1         4
290       1         4
291       1         4
292       1         4
293       1         4
294       1         4
295       1         4
296       1         4
297       1         4
298       1         4
299       1         4
300       1         4
301       1         4
302       1         4
303       1         4
304       1         4
305       1         4
306       1         4
307       1         4
308       1         4
309       1         4
310       1         4
311       1         4
312       1         4
313       1         4
314       1         4
315       1         4
316       1         4
317       1         4
318       1         4
319       1         4
320       1         4
321       1         4
322       1         4
323       1         4
324       1         4
325       1         4
326       1         4
327       1         4
328       1         4
329       1         4
330       1         4
331       1         4
332       1         4
333       1         4
334       1         4
335       1         4
336       1         4
337       1         4
338       1         4
339       1         4
340       1         4
341       1         4
342       1         4
343       1         4
344       1         4
345       1         4
346       1         4
347       1         4
348       1         4
349       1         4
350       1         4
351       1         4
352       1         4
353       1         4
354       1         4
355       1         4
356       1         4
357       1         4
358       1         4
359       1         4
360       1         4
361       1         4
362       1         4
363       1         4
364       1         4
365       1         4
366       1         4
367       1         4
368       1         4
369       1         4
370       1         4
371       1         4
372       1         4
373       1         4
374       1         4
375       1         4
376       1         4
377       1         4
378       1         4
379       1         4
380       1         4
381       1         4
382       1         4
383       1         4
384       1         4
385       1         4
386       1         4
387       1         4
388       1         4
389       1         4
390       1         4
391       1         4
392       1         4
393       1         4
394       1         4
395       1         4
396       1         4
397       1         4
398       1         4
399       1         4
400       1         4

```

• Data Preprocessing Done

Data Preprocessing is the major step before model building as without proper cleaning ,proper formatting and understanding the data we cannot do good analysis. So Data Preprocessing being the major step is done by following ways :

- Importing the suitable libraires and then load the data
- Checking null values
- Dropping duplicated data

Importing the suitable libraires and then load the data:

First step in data preprocessing is importing the libraires which will be helpful in analysis.After loading we will check the shape ,information of dataset.

Importing Libraries

```

In [1]: 1 import pandas as pd
          2 import numpy as np
          3 import matplotlib.pyplot as plt
          4 %matplotlib inline
          5 import seaborn as sns
          6 import string
          7 import re
          8 from nltk.corpus import stopwords
          9 from nltk.tokenize import regexp_tokenize
          10 from wordcloud import WordCloud
          11
          12 from nltk.stem import WordNetLemmatizer
          13 import warnings
          14 warnings.filterwarnings('ignore')
          15
          16 # packages from gensim
          17 from gensim import corpora
          18 from gensim.parsing.preprocessing import STOPWORDS
          19

```

Loading the data

Loading dataset

```
In [2]: 1 # Loading smartphones dataset
2 smartphones = pd.read_csv(r"C:\Users\siwac\machine learning models\smartphones.csv")
3 smartphones.head(5)
```

```
Out[2]:
```

	Reviews	Ratings
0	Please be aware that this is a downgraded vers...	2
1	What a Smartphone by Xiaomi! I am thoroughly l...	5
2	Received my Mi 11x Pro today 🥳. Really excited...	4
3	I love gaming and have been waiting for a SDM ...	5
4	I am a oneplus 5 user for the last 4 years. So...	3

```
In [3]: 1 # Loading laptops dataset
2 laptops = pd.read_csv(r"C:\Users\siwac\machine learning models\laptops.csv")
3 laptops.head(5)
```

```
Out[3]:
```

	Reviews	Ratings
0	The laptop came in excellent package & the pac...	5
1	Its superb its world's lightest laptopSuperb ba...	5
2	Excellent W & B. Good design, good for the. P...	5

Checking null values :

As now we have loaded the data we will check the null values present in dataset. For checking the null values we will use heatmap.

```
In [8]: 1 # Concatenating all dataset together
2 df = pd.concat([smartphones,laptops,iphone,refrigerators,monitors,headphones],ignore_index=True)
3 df.shape
```

```
Out[8]: (39752, 2)
```

```
In [9]: 1 # Removing all duplicated data from dataset
2 df.drop_duplicates(inplace=True,ignore_index=True)
```

```
In [10]: 1 # Checking null values
2 df.isnull().sum()
```

```
Out[10]: Reviews    5
Ratings    0
dtype: int64
```

```
In [11]: 1 # Dropping null values
2 df.dropna(inplace=True)
```

```
In [12]: 1 # Checking data info
2 df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 18516 entries, 0 to 18510
Data columns (total 2 columns):
 #   Column  Non-Null Count  Dtype

```

So there are no null values in our dataset.

• Hardware and Software Requirements and Tools Used

There is no such requirement for hardware ,but I have used intel i5 8th generation processor.

Software: Jupyter Notebook (Anaconda 3)

Language : Python 3.9

Libraries used in project:

- a. Pandas
- b. Gensim
- c. Nltk
- d. Seaborn
- e. Sklearn
- f. Re

• Data Cleaning

- As we are working on text data so NLP technique_ is used here to deal with the data. Using this technique we will clean data into simple steps.
- STEP 1 : Converting textual data into lower case letters.
- STEP 2:Removal of stopwords and punctuations
- STEP 3: Converting text into tokens using a method `regexp_tokenize`
- STEP 4 Lemmatization(converting words into their simpler forms)

Converting all upper case letters into lower case letters:

- We have used `str.lower()` method to convert the uppercase letters into lower case letters.

```
In [19]: 1 # Lowering the upper case letters
          2
          3 reviews['Reviews'] = reviews['Reviews'].str.lower()
```

```
In [20]: 1 # Checking sample review
          2 reviews['Reviews'][0]
```

```
Out[20]: 'please be aware that this is a downgraded version of globally launched mi 11 and mi 11 pro.reason for low rating -1. ordinary
display. this is not the 2k display as featured on global versions of mi 11 series1. no gorilla glass victus unlike global vers
ion.3. band tests shows poor build quality which may result in breakage of phone due to bending4. such a cheap tactic of sellin
g a phone with the promise of fast charging and not providing charger with a adequate charging capacity.5. no curved display un
like global version.6. took almost 20 days to get my hands on the phone post order.7. reported heating issues in sd 888 while p
laying intensive gamesin short, xiaomi has attempted to sell big story by first launching mi 11 globally and then putting "x" a
nd selling subpar product in india through illusion.i guess, xiaomi should not have cut the corners and compromise the quality
in indian versions. indians are willing to pay premium and success of oneplus is the testament to that. any cunning move to und
ercut the value proposition will hurt xiaomi in long run.(note that i currently use redmi note 10 pro and have used so many pho
nes of xiaomi ever since it launched first ever phone in india and hence hold a adequate credentials to voice my criticism as a
2022/11/11 11:11:11'
```


Removal of stopwords and then converting sentences into tokens

- After applying the lower case method we defined a function which will remove stopwords and convert the sentence data into tokens so that unnecessary text can be removed easily. Using the `regex_tokenize()` method we simply keep the data from a-z case letters and remove all numbers.

```
In [21]: 1 # Removal of stopwords and unnecessary words like punctuations, signs, numbers
2 def stop_words(text):
3     stop_wrds=set(STOPWORDS).union(set(string.punctuation))
4     tokens=regex_tokenize(text,"[a-z']+")
5     return [t for t in tokens if t not in stop_wrds]
```

```
In [22]: 1 # Apply the above function to the columns reviews
2 reviews['Reviews']=reviews["Reviews"].apply(stop_words)
```

```
In [23]: 1 reviews
```

```
Out[23]:
```

	Reviews	Ratings
0	{aware, downgraded, version, globally, launche...	2.0
1	{smartphone, xiaomi, thoroughly, impressed, pl...	5.0
2	{received, mi, x, pro, today, excited, came, p...	4.0
3	{love, gaming, waiting, sdm, k, compromise, as...	5.0
4	{oneplus, user, years, compare, oneplus, mi, x...	3.0

Lemmatizing the tokens (convert them into simpler forms)

- Lemmatization technique is used to convert the word into its simpler form like from plural to singular, from 3rd form of verb to 1st form of verb. We applied all the techniques and append the result into list which will be stored into a new column.

```
In [25]: 1 # After stopwords removal now we are going to lemmatize words
2
3 result = [] # empty list to store data
4 lemma = WordNetLemmatizer()
5
6 # function for lemmatization
7 def lemmatize(text):
8     word1=lemma.lemmatize(text,pos='n')
9     word2=lemma.lemmatize(word1,pos='v')
10    return lemma.lemmatize(word2,pos=('a'))
11
12 # function for filtering and apply above function to the column
13 def preprocess(text):
14     result=[]
15     for token in text:
16         if len(token)>=3:
17             result.append(lemmatize(token))
18
19     return result
```

Cleaned text

- After the completion of all the steps we are left with clean data. Now this data can be used for model building after converting textual data into vector form.

```
In [27]: 1 # Assigning new column to the processed reviews
2 reviews['cleaned_reviews'] = processed_review
3

In [28]: 1 # Changing uncleaned reviews into cleaned reviews
2 reviews['Reviews'] = reviews['cleaned_reviews'].apply(lambda x: ' '.join(y for y in x))

In [29]: 1 # checking sample of cleaned review
2 reviews['Reviews'][0]
```

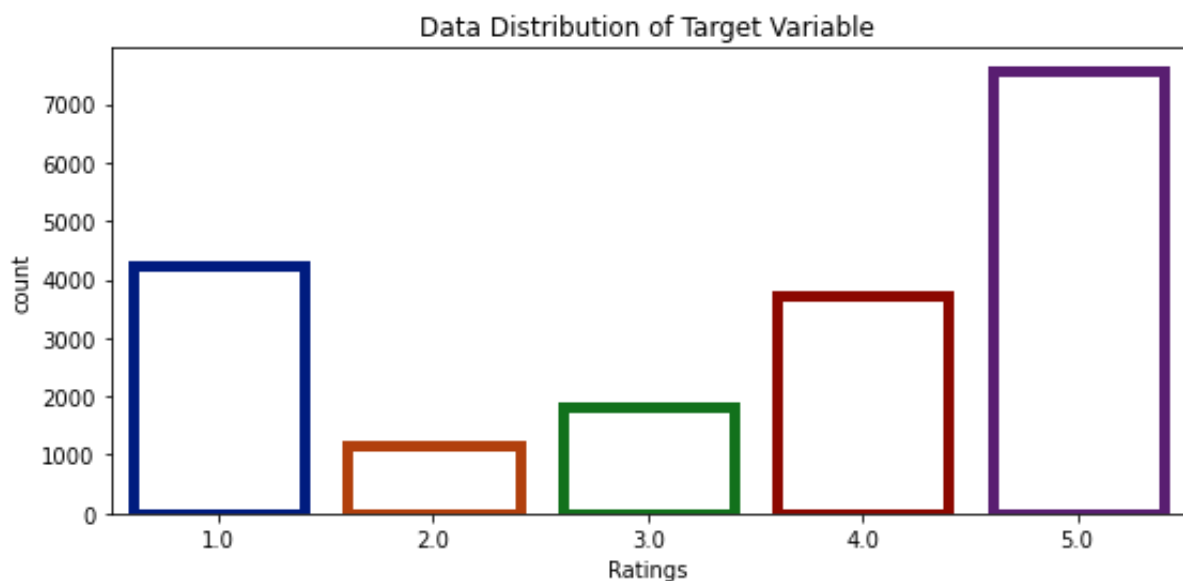
Out[29]: 'aware downgrade version globally launch pro reason low rat ordinary display display feature global version series gorilla glas
s victus unlike global version band test show poor build quality result breakage phone bend cheap tactic sell phone promise fas
t charge provide charger adequate charge capacity curve display unlike global version take day hand phone post order report hea
t issue play intensive gamesin short xiaomi attempt sell big story launch globally put sell subpar product india illusion guess
xiaomi cut corner compromise quality indian version indian will pay premium success oneplus testament cunning undercut value pr
oposition hurt xiaomi long run note currently use redmi note pro phone xiaomi launch phone india hold adequate credential voice
criticism fan'

• Exploratory Data Analysis

```
In [16]: 1 # Let's check total count of each rating
2 reviews['Ratings'].value_counts()

Out[16]: 5.0    7582
1.0    4251
4.0    3719
3.0    1816
2.0    1148
Name: Ratings, dtype: int64

In [17]: 1 # Let's check data distribution of target variable
2 plt.figure(figsize=(9,4),frameon=False)
3 sns.countplot('Ratings',data=reviews,color="salmon", facecolor=(0, 0, 0, 0),
4               linewidth=5,
5               edgecolor=sns.color_palette("dark", 5))
6 plt.title("Data Distribution of Target Variable")
7 plt.xlabel("Ratings")
8 plt.show()
```



Wordcloud to see the words which are most repeated in 5star ratings:

```
1 # Using wordcloud see the texts which are most repeated in 5 star ratings
2 star_5 = reviews['Reviews'][reviews['Ratings']==5.0]
3 combined_5star=' '.join(star_5)
4 plt.figure(figsize=(18,12),facecolor='grey')
5 wordcloud = WordCloud(max_font_size=50,max_words=250,background_color='white').generate(combined_5star)
6
7 # Display the generated image:
8 plt.imshow(wordcloud, interpolation='lanczos')
9 plt.axis("off")
10 plt.show()
```



Wordcloud to see the words which are most repeated in 4star ratings:

```
1 # Using wordcloud see the texts which are most repeated in 4 star ratings
2 star_4 = reviews['Reviews'][reviews['Ratings']==4.0]
3 combined_4star = ' '.join(star_4)
4 plt.figure(figsize=(18,12),facecolor='gray')
5 wordcloud = WordCloud(max_font_size=50,max_words=250,background_color='white').generate(combined_4star)
6
7 # Display the generated image:
8 plt.imshow(wordcloud, interpolation='lanczos')
9 plt.axis("off")
10 plt.show()
```



Wordcloud to see the words which are most repeated in 3star ratings:


```
1 # Using wordcloud see the texts which are most repeated in 3 star ratings
2 star_3 = reviews['Reviews'][reviews['Ratings']==3.0]
3 combined_3star=""'.join(star_3)
4 plt.figure(figsize=(18,12),facecolor='grey')
5 wordcloud = WordCloud(max_font_size=50,max_words=250,background_color='white').generate(combined_3star)
6
7 # Display the generated image:
8 plt.imshow(wordcloud, interpolation='lanczos')
9 plt.axis("off")
10 plt.show()
```



Wordcloud to see the words which are most repeated in 2star ratings:

```
1 # Using wordcloud see the texts which are most repeated in 2 star ratings
2 star_2 = reviews['Reviews'][reviews['Ratings']==2.0]
3 combined_2star=" ".join(star_2)
4 plt.figure(figsize=(18,12),facecolor='grey')
5 wordcloud = WordCloud(max_font_size=50,max_words=250,background_color='white').generate(combined_2star)
6
7 # Display the generated image:
8 plt.imshow(wordcloud, interpolation='lanczos')
9 plt.axis("off")
10 plt.show()
```



Wordcloud to see the words which are most repeated in 1star ratings:


```

In [38]: 1 # copying dataset
          2 df_model = reviews.copy()

In [39]: 1 # Importing library to vectorize the data for model building
          2 from sklearn.feature_extraction.text import TfidfVectorizer

In [40]: 1 # Applying function for vectorizing
          2 tf_idf1 = TfidfVectorizer(max_features=750,norm='l2')
          3

In [41]: 1 # fit and transform the dataset
          2 features = tf_idf1.fit_transform(df_model['Reviews'])
          3 features

Out[41]: <18516x750 sparse matrix of type '<class 'numpy.float64'>'
          with 323487 stored elements in Compressed Sparse Row format>

```

```

In [42]: 1 # storing in dataframe
          2 vectorized_features= pd.DataFrame(features.todense(),columns=tf_idf1.get_feature_names())
          3 vectorized_features

```

```

Out[42]:
      able absolutely  accord  accurate  acer  actually  adapter  add  adjust  adjustable  ...  worry  worth  write  wrong  xiaomi  year  yes  you  youth
0  0.0         0.0      0.0      0.0  0.0      0.0  0.00000  0.0   0.0      0.0  _  0.0   0.0   0.0   0.0  0.389902  0.000000  0.0  0.0
1  0.0         0.0      0.0      0.0  0.0      0.0  0.00000  0.0   0.0      0.0  _  0.0   0.0   0.0   0.0  0.173649  0.000000  0.0  0.0
2  0.0         0.0      0.0      0.0  0.0      0.0  0.21108  0.0   0.0      0.0  _  0.0   0.0   0.0   0.0  0.000000  0.000000  0.0  0.0
3  0.0         0.0      0.0      0.0  0.0      0.0  0.00000  0.0   0.0      0.0  _  0.0   0.0   0.0   0.0  0.121831  0.092367  0.0  0.0
4  0.0         0.0      0.0      0.0  0.0      0.0  0.00000  0.0   0.0      0.0  _  0.0   0.0   0.0   0.0  0.000000  0.209272  0.0  0.0
...
18511  0.0         0.0      0.0      0.0  0.0      0.0  0.00000  0.0   0.0      0.0  _  0.0   0.0   0.0   0.0  0.000000  0.000000  0.0  0.0
18512  0.0         0.0      0.0      0.0  0.0      0.0  0.00000  0.0   0.0      0.0  _  0.0   0.0   0.0   0.0  0.000000  0.000000  0.0  0.0
18513  0.0         0.0      0.0      0.0  0.0      0.0  0.00000  0.0   0.0      0.0  _  0.0   0.0   0.0   0.0  0.000000  0.000000  0.0  0.0
18514  0.0         0.0      0.0      0.0  0.0      0.0  0.00000  0.0   0.0      0.0  _  0.0   0.0   0.0   0.0  0.000000  0.000000  0.0  0.0
18515  0.0         0.0      0.0      0.0  0.0      0.0  0.00000  0.0   0.0      0.0  _  0.0   0.0   0.0   0.0  0.000000  0.000000  0.0  0.0

18516 rows x 750 columns

```

Feature selection :


```

In [47]: 1 # Using feature selection for better model building
          2 selection = SelectFromModel(LogisticRegression())
          3 selection.fit(vectorized_features,y)

Out[47]: SelectFromModel(estimator=LogisticRegression())

In [48]: 1 # finding valuable features
          2 feature_name = vectorized_features.columns[selection.get_support()]
          3 feature_name

Out[48]: Index(['absolutely', 'adapter', 'add', 'amaze', 'amazon', 'amd', 'amoled',
               'anc', 'android', 'angle',
               ...,
               'weight', 'wide', 'wifi', 'wire', 'wireless', 'wise', 'worry', 'wrong',
               'xiaomi', 'yes'],
              dtype='object', length=297)

```

```

In [49]: 1 # Apply transformation on train independent dataset.
          2
          3 selection_New_x = selection.transform(vectorized_features)
          4 New_x = pd.DataFrame(selection_New_x,columns=feature_name)
          5 New_x

Out[49]:

```

	absolutely	adapter	add	amaze	amazon	amd	amoled	anc	android	angle	...	weight	wide	wifi	wire	wireless	wise	worry	wrong	x
0	0.0	0.000000	0.0	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	...	0.000000	0.0	0.0	0.000000	0.000000	0.0	0.0	0.0	0.38
1	0.0	0.000000	0.0	0.120381	0.0	0.0	0.162607	0.0	0.0	0.0	...	0.000000	0.0	0.0	0.000000	0.000000	0.0	0.0	0.0	0.11
2	0.0	0.21108	0.0	0.000000	0.0	0.0	0.195027	0.0	0.0	0.0	...	0.000000	0.0	0.0	0.000000	0.000000	0.0	0.0	0.0	0.08
3	0.0	0.000000	0.0	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	...	0.000000	0.0	0.0	0.000000	0.000000	0.0	0.0	0.0	0.11
4	0.0	0.000000	0.0	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	...	0.000000	0.0	0.0	0.000000	0.000000	0.0	0.0	0.0	0.08
...
18511	0.0	0.000000	0.0	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	...	0.000000	0.0	0.0	0.000000	0.157517	0.0	0.0	0.0	0.08
18512	0.0	0.000000	0.0	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	...	0.000000	0.0	0.0	0.000000	0.000000	0.0	0.0	0.0	0.08
18513	0.0	0.000000	0.0	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	...	0.000000	0.0	0.0	0.000000	0.000000	0.0	0.0	0.0	0.08
18514	0.0	0.000000	0.0	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	...	0.099687	0.0	0.0	0.308569	0.115878	0.0	0.0	0.0	0.08
18515	0.0	0.000000	0.0	0.000000	0.0	0.0	0.000000	0.0	0.0	0.0	...	0.000000	0.0	0.0	0.000000	0.000000	0.0	0.0	0.0	0.08

18516 rows x 297 columns

Splitting the dataset into training and testing

```
In [50]: 1 # Let's now find the best random state
2 max_f1_score=0
3 for i in range(1,100):
4     x_train,x_test,y_train,y_test=train_test_split(New_x,y,test_size=.22,random_state=i)
5     mnb=MultinomialNB()
6     mnb.fit(x_train,y_train)
7     y_pred=mnb.predict(x_test)
8     f1=f1_score(y_test,y_pred, average='weighted')
9
10     if f1>max_f1_score:
11         max_f1_score=f1
12         final_random_state=i
13 print("Max f1 score corresponding to ",final_random_state,"is ",max_f1_score,".")
Max f1 score corresponding to 89 is 0.4289328029730464 .

In [51]: 1 # Let's split the train and test dataset
2 x_train,x_test,y_train,y_test = train_test_split(New_x,y,test_size=0.2,random_state=final_random_state)
3 print("Train shapes : X = {}, y = {}".format(x_train.shape,y_train.shape))
4 print("Test shapes : X = {}, y = {}".format(x_test.shape,y_test.shape))
5 # Success
6 print("\nTraining and testing split was successful.")
Train shapes : X = (14812, 297), y = (14812,)
Test shapes : X = (3704, 297), y = (3704,)

Training and testing split was successful.
```

• Model Evaluation

```
In [52]: 1 # function for finding model's accuracy and other metrics
2
3 def metrics_score(model):
4     print(model)
5     model.fit(x_train,y_train)
6     print("\nTraining ",round(model.score(x_train,y_train)*100,1),"%")
7     y_pred=model.predict(x_test)
8     print("\nAccuracy score - ",round(accuracy_score(y_test,y_pred)*100,1),"%")
9     print("\nClassification Report - \n{}".format(classification_report(y_test,y_pred)))
10
11     print("\n\nCross Validation Score - ")
12     for i in range(2,10):
13         cross_val=cross_val_score(model,New_x,y,cv=i,n_jobs=1) # cross validation
14
15         print("Cross validation score at ",i,"is",round(cross_val.mean()*100,1),"%")
16         print(round(cross_val.std()*100,1))
17
```

```
In [53]: 1 metrics_score(MultinomialNB(alpha=1.0,fit_prior=True))
MultinomialNB()
Training 52.3 %
Accuracy score - 53.7 %
Classification Report -
      precision    recall  f1-score   support

      1.0         0.64      0.58      0.61         893
      2.0         0.00      0.00      0.00         200
      3.0         1.00      0.00      0.01         356
      4.0         0.27      0.02      0.03         712
      5.0         0.51      0.95      0.66        1543

 accuracy          0.54        3704
 macro avg          0.48        3704
 weighted avg       0.52        3704
```


GaussianNB()

Training 45.7 %

Accuracy score - 43.3 %

Classification Report -

	precision	recall	f1-score	support
1.0	0.48	0.58	0.52	893
2.0	0.12	0.26	0.17	200
3.0	0.18	0.26	0.21	356
4.0	0.28	0.16	0.21	712
5.0	0.64	0.54	0.59	1543
accuracy			0.43	3704
macro avg	0.34	0.36	0.34	3704
weighted avg	0.46	0.43	0.44	3704

Cross Validation Score -

Cross validation score at 2 is 22.8 %

0.4

Cross validation score at 3 is 32.7 %

4.9

Cross validation score at 4 is 34.2 %

BernoulliNB()

Training 52.2 %

Accuracy score - 51.5 %

Classification Report -

	precision	recall	f1-score	support
1.0	0.52	0.71	0.60	893
2.0	0.22	0.11	0.15	200
3.0	0.36	0.15	0.21	356
4.0	0.30	0.20	0.24	712
5.0	0.60	0.69	0.64	1543
accuracy			0.51	3704
macro avg	0.40	0.37	0.37	3704
weighted avg	0.48	0.51	0.48	3704

Cross Validation Score -

Cross validation score at 2 is 47.1 %

4.5

Cross validation score at 3 is 47.3 %

4.7

Cross validation score at 4 is 48.7 %

6.1

LogisticRegression()

Training 55.3 %

Accuracy score - 55.0 %

Classification Report -

	precision	recall	f1-score	support
1.0	0.60	0.65	0.62	893
2.0	0.08	0.01	0.01	200
3.0	0.37	0.07	0.12	356
4.0	0.37	0.15	0.21	712
5.0	0.56	0.86	0.68	1543
accuracy			0.55	3704
macro avg	0.39	0.35	0.33	3704
weighted avg	0.49	0.55	0.49	3704

Cross Validation Score -

Cross validation score at 2 is 51.5 %

3.3

Cross validation score at 3 is 52.8 %

4.0

Cross validation score at 4 is 52.5 %

- -

Hyperparameter Tunning :

```
In [59]: 1 # Defining parameters
2 solvers = ['newton-cg', 'lbfgs', 'liblinear']
3 penalty = ['l2']
4 c_values = [100, 10, 1.0, 0.1, 0.01]

In [60]: 1 from sklearn.model_selection import RepeatedStratifiedKFold, GridSearchCV
2
3 # define grid search
4 grid = dict(solver=solvers,penalty=penalty,C=c_values)
5 # Cross validation technique
6 cv_method = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
7
8 grid_search = GridSearchCV(estimator=LogisticRegression(), param_grid=grid, n_jobs=-1, cv=cv_method, scoring='accuracy', error_score='raise')
9 grid_result = grid_search.fit(New_x, y)

In [61]: 1 # summarize results
2 print("Best: %f using %s" % (grid_result.best_score_, grid_result.best_params_))
3 means = grid_result.cv_results_['mean_test_score']
4 stds = grid_result.cv_results_['std_test_score']
5 params = grid_result.cv_results_['params']
6 for mean, stdev, param in zip(means, stds, params):
7     print("%f (%f) with: %r" % (mean, stdev, param))
```

```
Best: 0.538508 using {'C': 1.0, 'penalty': 'l2', 'solver': 'liblinear'}
0.534187 (0.008788) with: {'C': 100, 'penalty': 'l2', 'solver': 'newton-cg'}
0.533359 (0.009112) with: {'C': 100, 'penalty': 'l2', 'solver': 'lbfgs'}
0.536924 (0.008259) with: {'C': 100, 'penalty': 'l2', 'solver': 'liblinear'}
0.535231 (0.008989) with: {'C': 10, 'penalty': 'l2', 'solver': 'newton-cg'}
0.535249 (0.009186) with: {'C': 10, 'penalty': 'l2', 'solver': 'lbfgs'}
0.537428 (0.008240) with: {'C': 10, 'penalty': 'l2', 'solver': 'liblinear'}
0.537770 (0.008779) with: {'C': 1.0, 'penalty': 'l2', 'solver': 'newton-cg'}
0.537950 (0.008973) with: {'C': 1.0, 'penalty': 'l2', 'solver': 'lbfgs'}
0.538508 (0.007772) with: {'C': 1.0, 'penalty': 'l2', 'solver': 'liblinear'}
0.528174 (0.006082) with: {'C': 0.1, 'penalty': 'l2', 'solver': 'newton-cg'}
0.528210 (0.006023) with: {'C': 0.1, 'penalty': 'l2', 'solver': 'lbfgs'}
0.524016 (0.005959) with: {'C': 0.1, 'penalty': 'l2', 'solver': 'liblinear'}
0.462465 (0.004260) with: {'C': 0.01, 'penalty': 'l2', 'solver': 'newton-cg'}
0.462465 (0.004260) with: {'C': 0.01, 'penalty': 'l2', 'solver': 'lbfgs'}
0.442176 (0.003898) with: {'C': 0.01, 'penalty': 'l2', 'solver': 'liblinear'}
```

```
1 # Using the best parameters obtained
2 model = LogisticRegression(C= 1.0, penalty='l2', solver='liblinear')
3 model.fit(x_train,y_train)
4 y_pred = model.predict(x_test)
5 print("Accuracy score = ",round(accuracy_score(y_test,y_pred)*100,1),"%")
6 print('Cross validation score: ',cross_val_score(model,New_x,y,cv=3,scoring='accuracy').mean()*100)
7 print("Confusion matrix\n",confusion_matrix(y_test,y_pred))
8 print("Classification report\n",classification_report(y_test,y_pred))
```

Accuracy score = 54.8 %

Cross validation score: 52.91099589544178

Confusion matrix

```
[[ 572    4    5   23  289]
 [ 100    1   14   21   64]
 [ 101    0   19   45  191]
 [  73    0    6   86  547]
 [ 107    0    7   79 1350]]
```

Classification report

	precision	recall	f1-score	support
1.0	0.60	0.64	0.62	893
2.0	0.20	0.01	0.01	200
3.0	0.37	0.05	0.09	356
4.0	0.34	0.12	0.18	712
5.0	0.55	0.87	0.68	1543
accuracy			0.55	3704
macro avg	0.41	0.34	0.32	3704
weighted avg	0.49	0.55	0.48	3704

CONCLUSION

- Key Findings and Conclusions of the Study

Working on this project was a wonderful experience as we have to collect data of reviews and ratings of customers who are using online retailer websites for shopping. After scrapping the data , we cleaned it by removing the stopwords, punctuations ,numbers and other unnecessary words that were not useful for our model building. While cleaning the dataset we find out that there are a lot of noise present in customer reviews.

Major problem we faced was in data cleaning because customers are not using proper language format or using other languages in reviews which

creates a major problem in identifying the real words and this leads to create more features.

- * We have collected a lot of data from e-commerce websites. Although for better model building it should be of different websites but due to the size and convenience i have used data from Amazon only.
- * After collecting all the data we combine it together into a dataframe and find out that duplicated data is present.
- * Then after using Gensim stopwords library we removed the unnecessary words.
- * After removal of stopwords we tokenize the sentence into words.
- * Then using the Lemmatization technique we converted words into thier simpler form.
- * After lemmatization we converted the textual data into vector form for model building.
- * Now we splited the dataset into training and testing dataset for model building.
- * Tried a lot of algorithms on dataset but `Bernoulli Naive Bayes`, `Multibnomail Naive Bayes` and `Logistic Regression` gives the better result.
- * We finalised the `Logistic Regression` model and then do hyperparmater tunning by `GridSearchCV` for findig best result out of it

- Suggestions:

- a. Less time complexity
- b. More accurate reviews can be given
- c. Less errors can be avoided.

Thankyou