

## **Query-based comparative Analysis of a Particular Product on Amazon and Flipkart & Sentiment Analysis on Desired Product**

### **Research Project Report**

**Shreyansh Padarha**

(21112026)

3 BSC DS A

Under the Guidance of  
Prof Naived George Eapen  
Department of Data Science

A Project Report submitted in Partial Fulfilment of the requirements for the award of  
Degree of Bachelor of Science (Data Science) of  
CHRIST (Deemed to be University), Pune Lavasa Campus - 'The Hub of Analytics'

2022 - 23



**CHRIST**  
(DEEMED TO BE UNIVERSITY)

PUNE LAVASA CAMPUS  
The Hub of Analytics

## CERTIFICATE

This is to certify that the research project titled Query-based Comparative Analysis of a Particular Product on Amazon and Flipkart & Sentiment Analysis on Desired Product is a Bonafede record of work done by Shreyansh Padarha (21112026) of 3 BSC DS, CHRIST (Deemed to be University), Pune Lavasa Campus, in partial fulfilment of the requirements of III Semester BSc (Data Science) during the year 2022-2023.

Dr Lija Jacob  
Head of the Department

Naived George Eapen  
Project Guide

## **Table of Contents**

1. Introduction
2. Background
3. Problem Statement
4. Data Sources and Collection
5. Tools and Libraries
6. Methodology
7. Observation and Findings
8. Details of Deployment and Presentation
9. Conclusion and Future Scope
10. References

## **[1] Introduction**

Customer confusion has been on the uprising for the past decade, this psychological construct can be credited to the rampant e-commerce industry. In 2021, e-retail sales crossed a mammoth 5.2 trillion USD worldwide. Although the numbers can be assumed to be bloated, as a customer, the plethora of e-commerce websites to choose from has made purchasing decisions harder. Ill-informed decisions regarding buying a product can lead to mistrust in online shopping or financial disasters. It's imperative to keep the customers well-informed to keep them satisfied, it also saves ample of time while purchasing a product.

## **[2] Background**

The three underlying principles of e-confusion or customer confusion are similarity confusion, overload confusion and unclarity confusion. But, in most developing countries like India, the biggest cause of consumer swimmers' is product price, rating, reviews and of course availability. Keeping this in mind, the project is aimed at building a thorough application to help users in deciding which product to buy from, and more importantly where from. The added feature of generating a sentiment analysis of the reviews made by customers, provides an in-depth understanding of the position of customers regarding the product. Opinion mining was an agenda that was exploited during the project, it helped in determining whether the overall customer experience of the customer is positive, neutral or negative in nature.

## **[3] Problem Statement**

This project attempts to combine data analytics methodologies like KDD (Knowledge Discovery Database) and comparative analysis techniques to tackle an ever-lasting conundrum, that is settling the question, which platform is the best to buy your desired product from? Considering time constraints, the comparisons were limited to Amazon and Walmart subsidiary Flipkart. The same is done by web scraping the requested products query from Amazon and Flipkart, and further web-scraping product reviews, enabling individual product-based sentimental analysis (Amazon Only).

## **[4] Data Sources and Collection**

The impromptu response nature of the project required an on-demand and customisable dataset generation, every time the user seeks a product. To tackle the complexity, we created three modules of web-scraping for specific purposes and generating specific datasets.

The sources of the web-scraped data were Amazon and Flipkart's website. The collection involved a generalised algorithm of web scrapping products and their information.

The reviews scrapped module, went through all the comments of the product and, post data scrubbing, storing the information in a dataset. Post scrapping the reviews the sentiment analysis module created by me, added 8 columns based on a ROBERTA and VADER model

sentiment scores. These scores helped in generating an aggregate opinion (polarity) representative of all the customers.

There were two main datasets that were created through the process.

- Web Scrapping Amazon (based on query (product given by client))
- Web Scrapping Flipkart (based on query (product given by client))
  - The above datasets, are merged and the final product dataset created with the web scrapped data had the variables:
    1. querry\_searched
    2. ecommerce\_website
    3. product\_name
    4. product\_price
    5. product\_rating
    6. rating\_count
    7. poduct\_image\_url
    8. product\_url
- Web Scrapping Comments (based on the product selected from the above dataset) (Amazon Only).
  - Post applying NLP models like ROBERTA and VADER on the comments scrapped, to generate sentiment scores. The comments scrapped dataset had the variables:
    1. customer\_name
    2. review\_date
    3. review\_title
    4. review\_ratings
    5. review\_content
    6. review\_size
    7. sentiment
    8. roberta\_sentiment
    9. vader\_sentiment
    10. roberta\_neg
    11. roberta\_neu
    12. roberta\_pos
    13. vader\_neg
    14. vader\_pos
    15. vader\_comp

## [5] Tools and Libraries

IDLE Used: Python Shell, PyCharm, Jupyter Notebook

Python modules and libraries used:

- BeautifulSoup
- requests
- pandas
- numpy
- scipy
- seaborn
- random
- matplotlib
- warnings
- time
- datetime
- pltoly
- regex
- re
- smtplib
- Ipython
- nltk.sentiment
- tqdm.notebook
- transformers
- scipy.special (softmax)
- tabulate
- urlib.request
- PIL
- Wordcount

## [6] Methodology

The process and methodology involved during the project were a robust and in-depth exploration of various components and fragments of KDD and data analytics.

- Understanding requirements to handle the Problem Statement.  
This first vital step made me realise the complexity of the project. Web-scraping using Python Libraries had to be learnt from scratch using tutorials and library documentation.
- Designing the architecture of the program  
From the onset, it was understood that the project required an elaborate backend and a sophisticated design to help integrate it with the finesse of the front end. This prompted me to create 3 separate modules and import it into the final file as .py finals. Using concepts from OOPs would result in a simplified and easy to traceback software.

The backend modules include :

- **webscraping.py** module:
  - Web Scrapping Amazon Class (based on query (product given by the client))
    - Web Scrapping Flipkart Class (based on query (product given by the client))
      - Dataset Created with the web scraped data had the variables:
        1. querry\_searched
        2. ecommerce\_website
        3. product\_name
        4. product\_price
        5. product\_rating
        6. rating\_count
        7. poduct\_image\_url
        8. product\_url
    - Web Scrapping Comments Class (based on the product selected from the above dataset) (Amazon Only)
      - Dataset Created with the web scraped data had the variables:
        1. customer\_name
        2. review\_date
        3. review\_title
        4. review\_ratings
        5. review\_content
  - **ProjectVisualisations.py** module:
    - Consists of various functions that help in plotting the amazon vs Flipkart data.
    - Using Plotly, seaborn, and matplotlib to come up with 17 plots for a website vs website analysis.
  - **ProductReviewAnalysis.py** module:
    - This file contained was the main file for performing a sentiment analysis upon reviews/comments on a product (currently limited to amazon)
      - + Using ROBERTA Model (NLP)
      - + Using VADER Model (NLP)
      - + A weighted 60:40 approach (Roberta : Vader) to obtain a final sentiment value.
      - + 18 Visually Appealing Plots

- Inspecting elements of static HTML pages of Amazon and Flipkart and Scrapping based on trial and error

The web scraping was done in tandem with the requests module and BeautifulSoup module, this involved requesting the search page based on the user's query. The page retrieved then had to be scouted for required information. This was done with basic knowledge of the inbuilt methods of the libraries being used and on a trial-and-error basis. It was a tedious process as it required checking where the elements were present on the actual page and comparing it in my script.

- Cleaning data and transforming data while it's being scrapped

A major difference between traditional EDAs and models created by me in the past and this project was, I had to clean and make the data apt while it was being scrapped and before inserting it into the dataset. This involved using regex and string functions and logic to make the data presentable and comprehensible.

- Visualising Amazon vs Flipkart product range for the product

Using Plotly, seaborn, and matplotlib to paint an effective picture of how the products vary on Amazon and Flipkart for a user's query. I created various plots using user-defined functions and recalled all of them under one visualising parent function.

- Scrapping All Comments based on the item selected

During, the initial bout of scraping all products from Amazon and Flipkart, the onclick() URL was also stored, in case of further analysis. Based on the product selected for sentiment analysis the URL is used for sending a request, and at the bottom of the product details page on Amazon, we have a "see all reviews" button. A request is sent to that button's onclick() URL, and a loop is used to change the URL page no. and all the reviews are scrapped along with the writer's name, date and review stars given. All of this is stored orderly in a panda's data frame.

- Sentiment Analysis using Vader and Roberta Model

The comments were analysed with the help of NLTK module and hugging face models. I used two different models (VADER and Roberta) and a weighted 60:40 approach to extract an accurate sentiment for the product.

VADER (Valence Aware Dictionary for Sentiment Reasoning) is a Model that is sensitive to polarity (positive/negative) and intensity of a emotion. For example words like 'love', 'enjoy', 'happy', 'like' all convey a positive sentiment. Also, VADER is intelligent enough to understand the basic context of these words, such as "did not love" as a negative statement. It also understands the emphasis of capitalization and punctuation, such as "ENJOY". The drawback of VADER is that it's optimised for social media, and it doesn't take context or sarcasm into consideration.

ROBERTA is a model that is an enhanced version of Google's BERT model. BERT makes use of a Transformer, an attention mechanism that learns contextual relations between words (or sub-words) in a text. In its vanilla form, Transformer includes two separate mechanisms an encoder that reads the text input and a decoder that produces a prediction for the task. ROBERTA has minor tweaks added to BERT. The model was trained with a hugging face 58 million tweets dataset, finetuned for natural language processing.

- Visualising Sentiment Analysis

The following maps, graphs and comparisons are used to display the results of the sentiment analysis:

- + Word Frequency Chart
- + Sentiment (Positive, Negative, Neutral) Count Plot, Pie Chart
- + Weighted Sentiment Spread (distplot)
- + Roberta vs Vader Model Sentiment Score comparison
- + Review count for the product (Bar graph, Pie Chart)
- + Length of the review vs Product Rating
- + In depth comparison of Vader and Roberta Scores

- Interesting Fact (Learning)

The project made use of 81 Methods (User-Defined Functions) in total, for efficient and sophisticated coding. The functions and blocks of code were easier to debug as well.

## [7] Observations and Findings

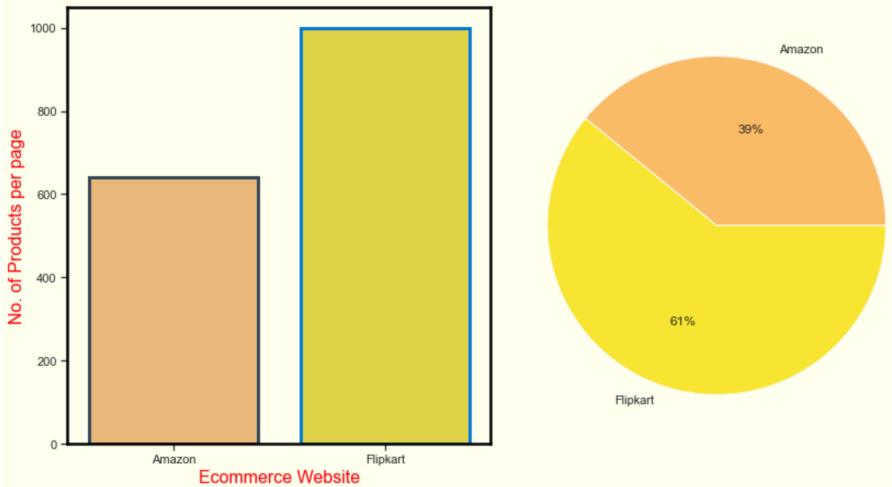
The observations are on a per-case basis, as can be seen in the below Classmate Notebook Comparison the results are self-explanatory.

-----AMAZON vs FLIPKART Comparison-----

PRODUCT QUERIED : Notebook  
BRAND : Classmate

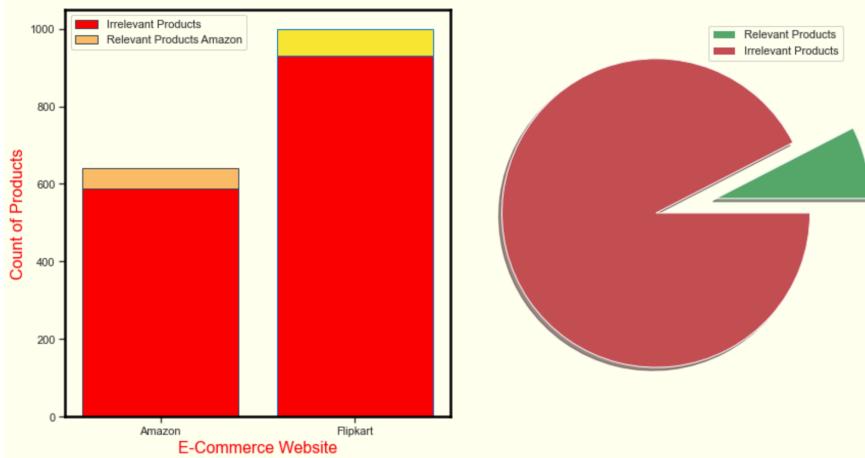
- > Amazon has 642 products listed on each page.
- > Flipkart has 1000 products listed on each page.

Products Listed Per Page for the Query

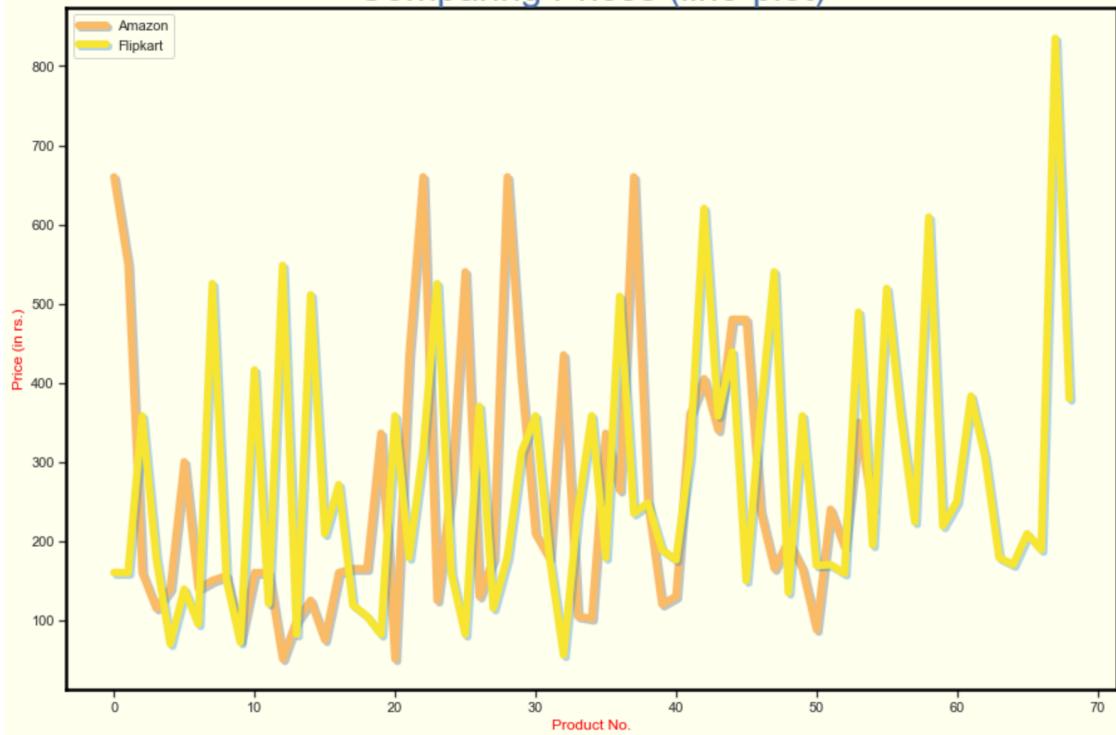


- > Amazon has 642 products listed on each page for the queried product, out of which 587 are of other brands
- > Flipkart has 1000 products listed on each page for the queried product, out of which 69 are of other brands

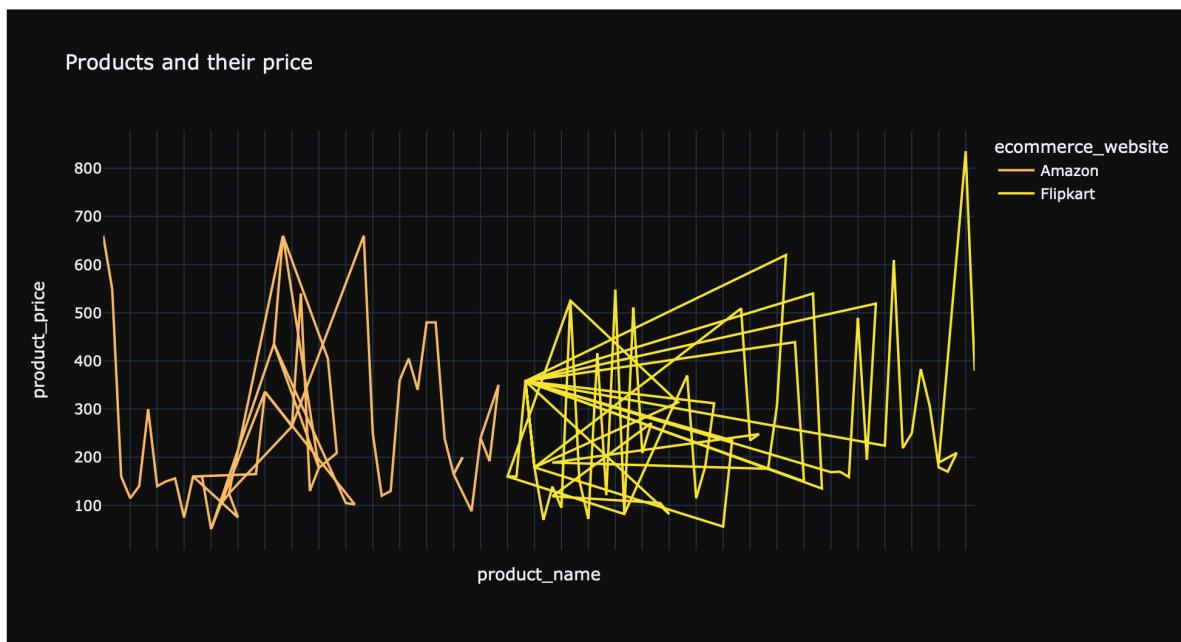
Stacked Bar Chart Comparing Relevance of Listed Products



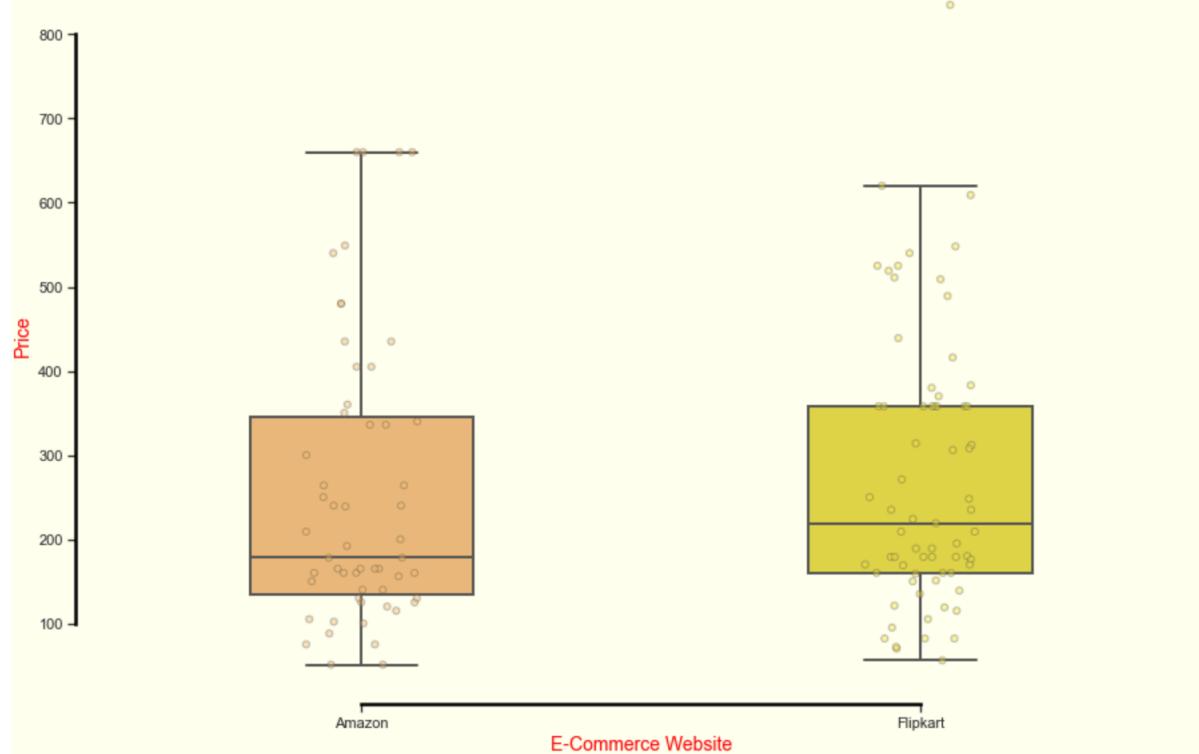
### Comparing Prices (line-plot)



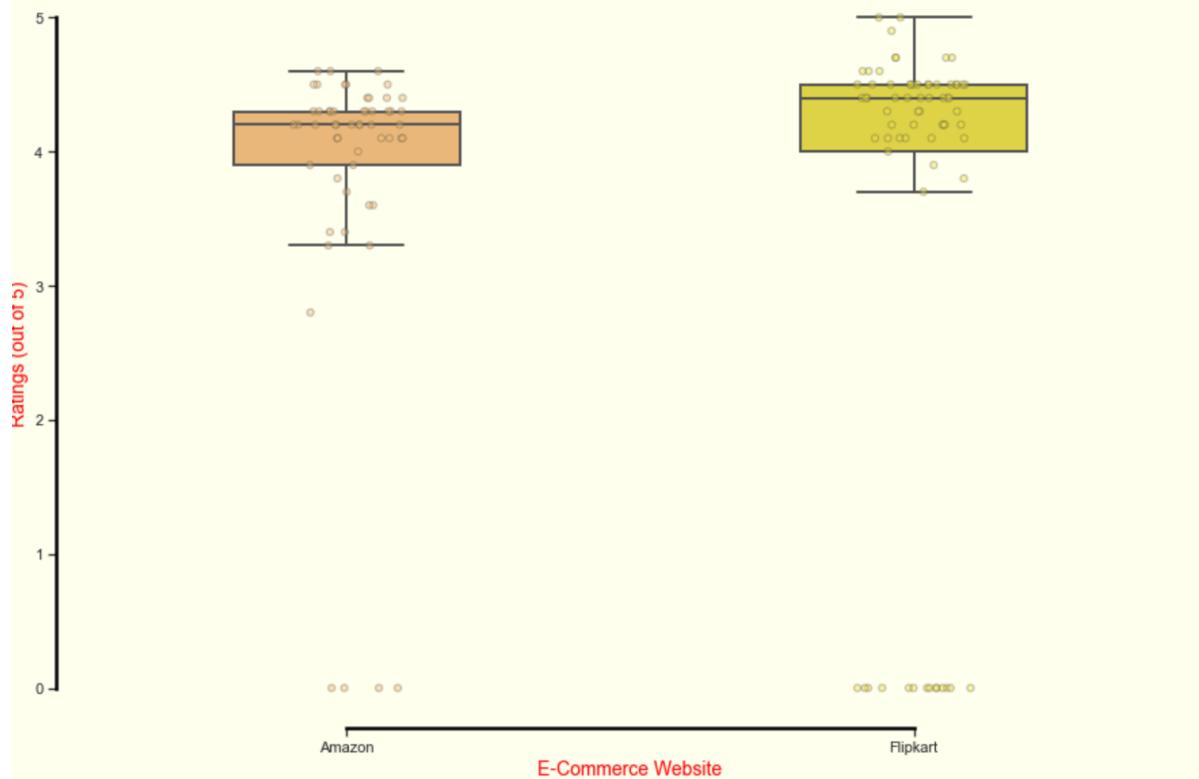
Average Product Price for queried product at amazon is : 256.7818181818182  
Average Product Price for queried product at flipkart is : 274.8985507246377



## Comparing Prices (box-plot)



## Comparing Ratings (box-plot)



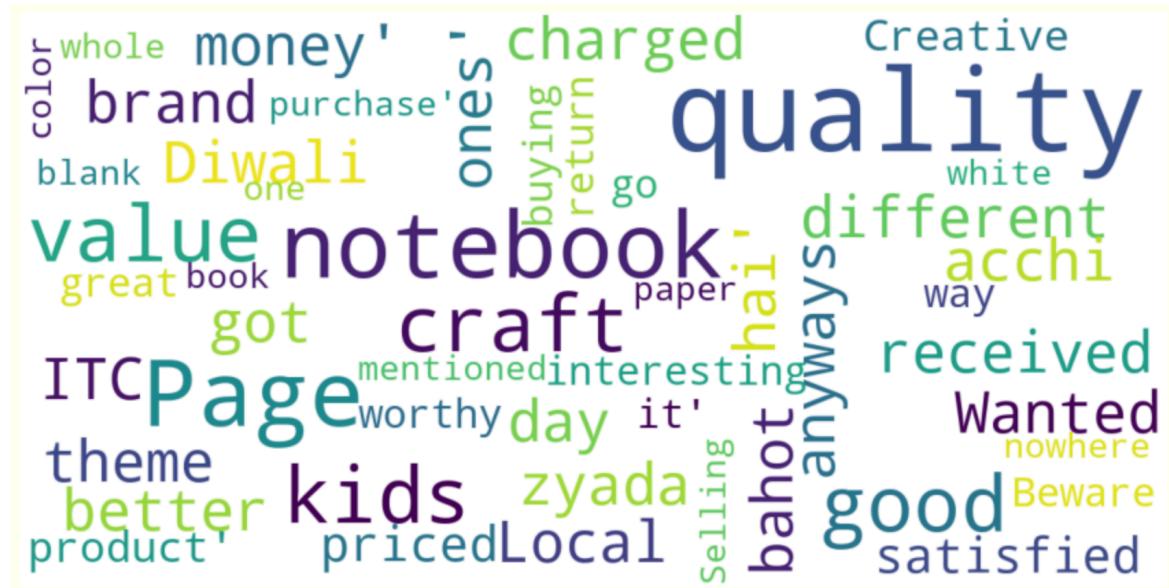
## Sentiment Analysis

### -----PRODUCT DETAILS-----

**Product Name :** Classmate Origami Notebooks – Single Line, 168 Pages, 240 mm x 180 mm – Pack Of 12  
**Product Price :** 660.0  
**Average Product Rating :** 4.3  
**Total Number of Ratings :** 0.0

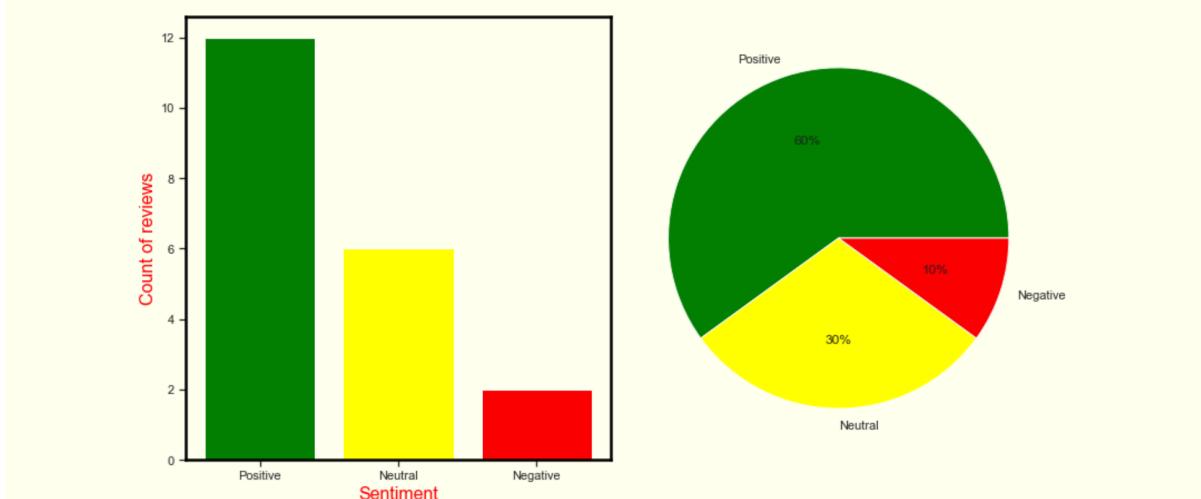


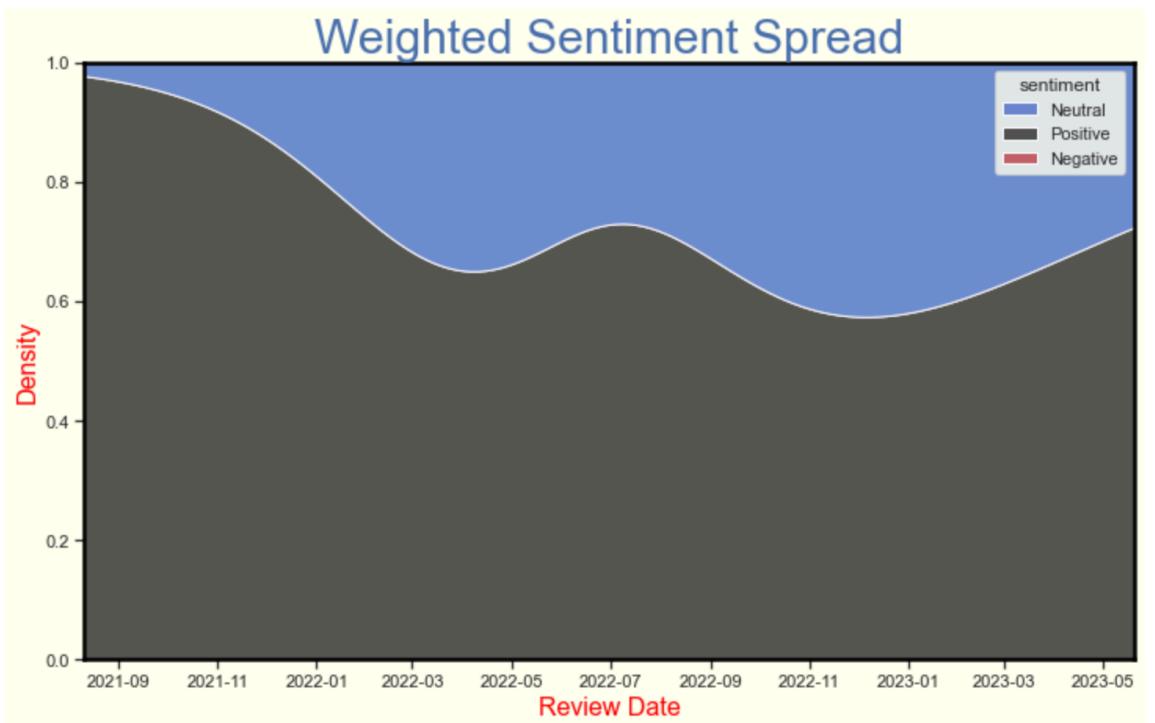
## WORD FREQUENCY CHART



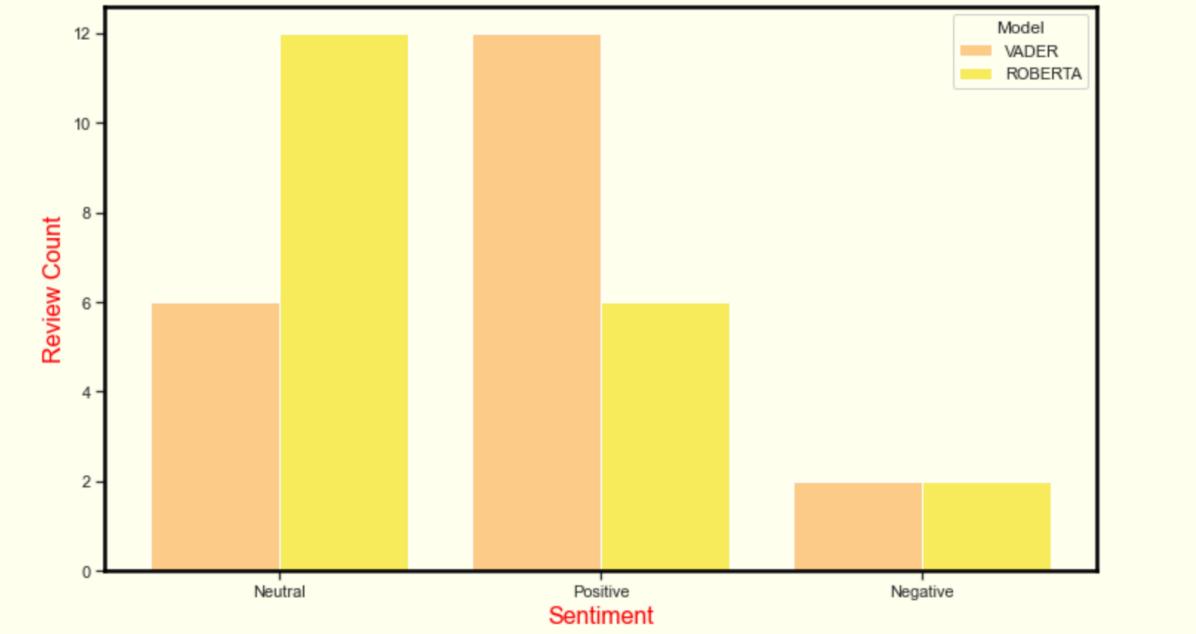
```
-> Positive Sentiment Reviews : 12  
-> Neutral Sentiment Reviews : 6  
-> Negative Sentiment Reviews : 2
```

Reviews Sentiment Chart (Weighted 3:2 (ROBERTA MODEL : VADER MODEL))

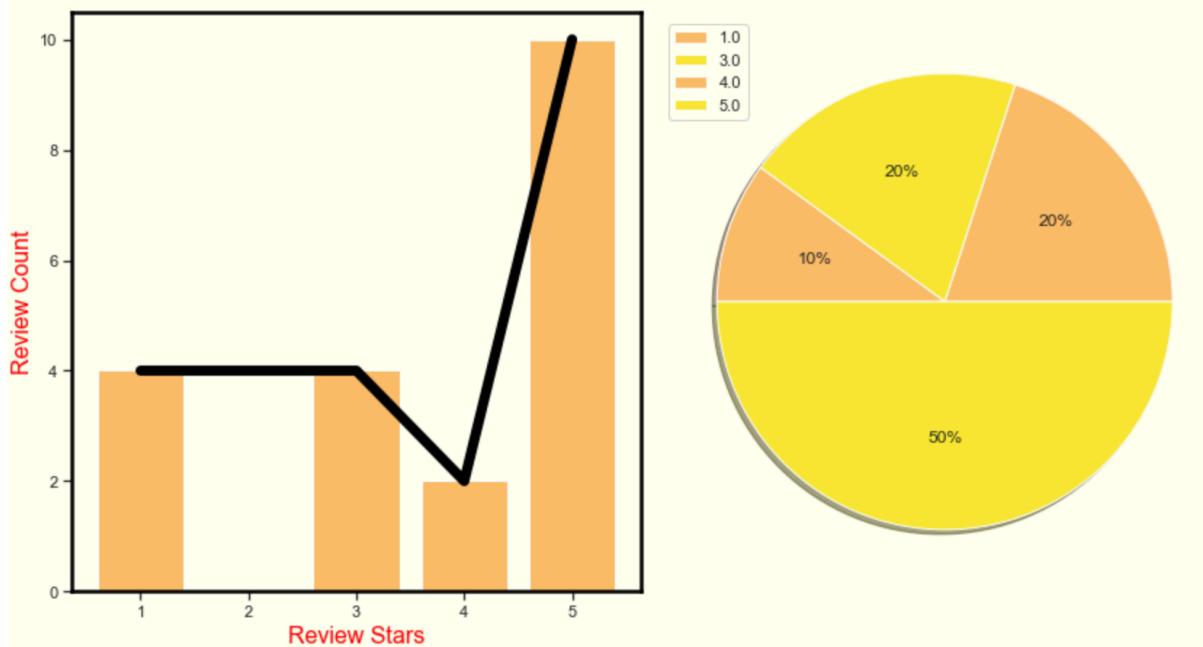




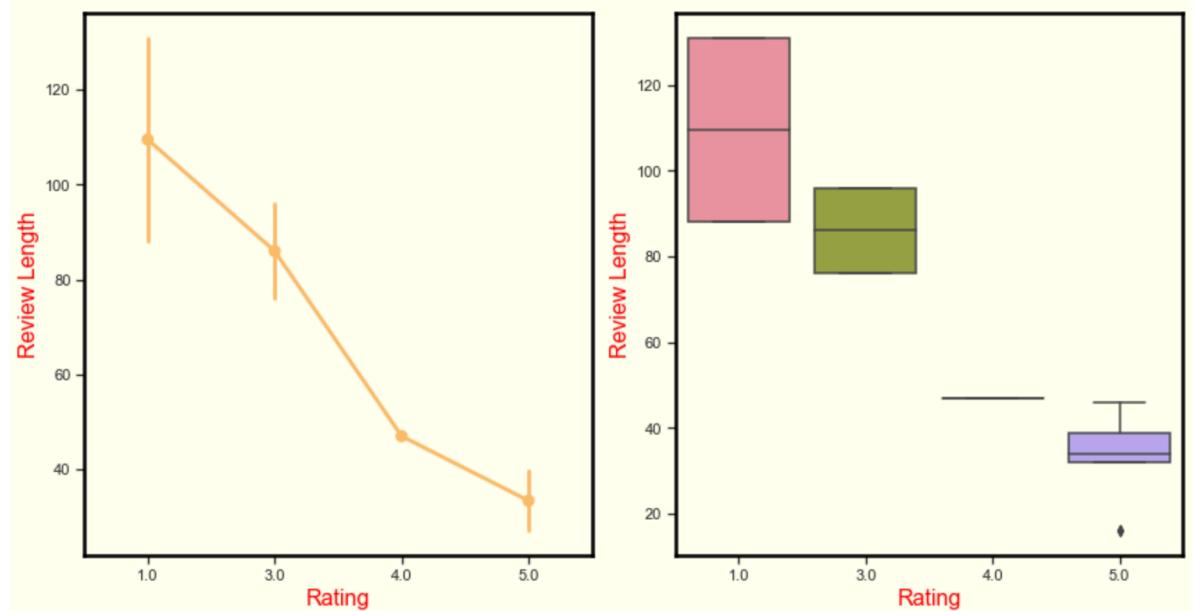
### ROBERTA VS VADER MODEL SENTIMENT COMPARISON



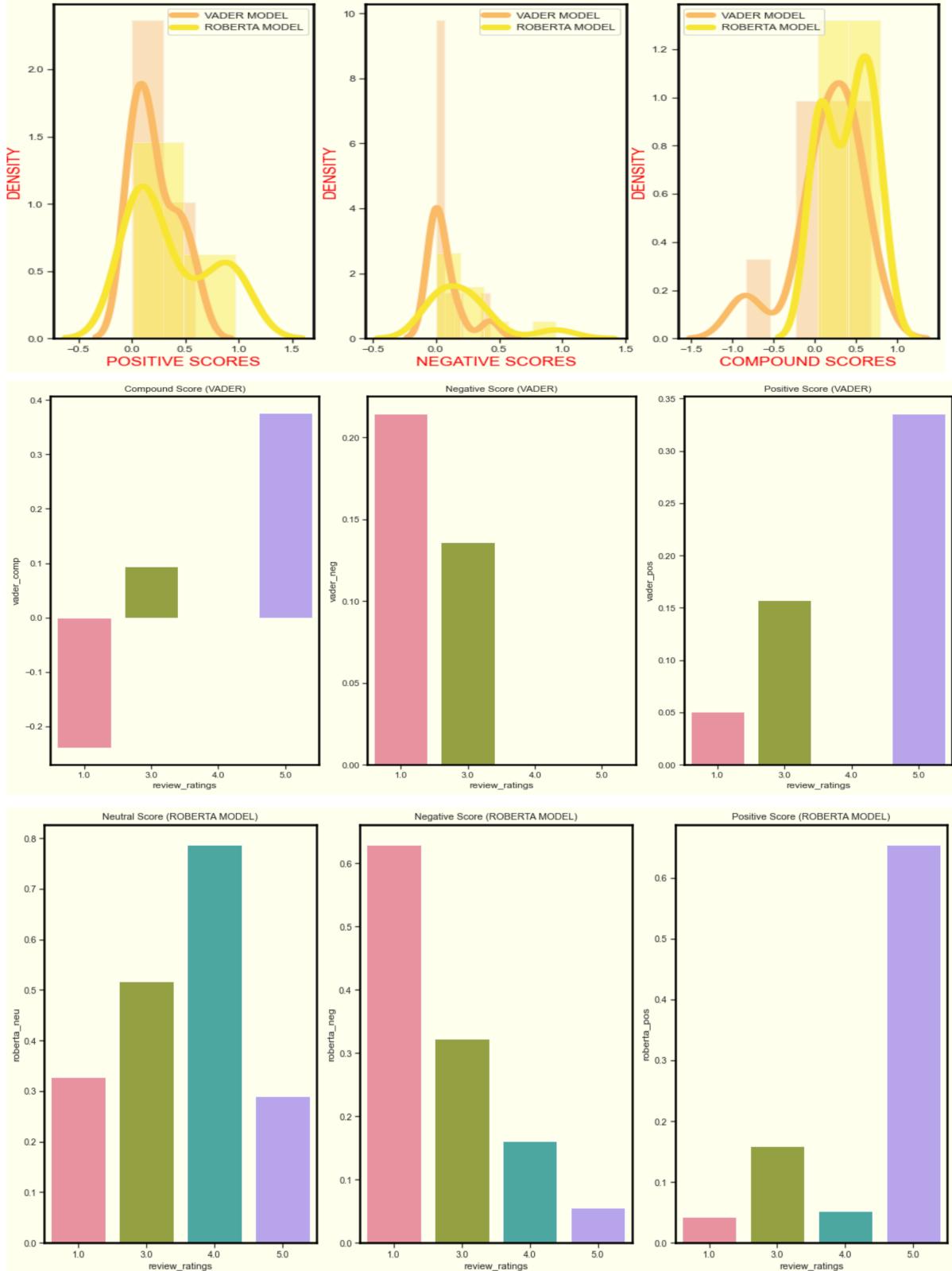
## Reviews Count for the product



## Product Rating vs Review Length



## In-Depth VADER and ROBERTA Model Analysis



## **[8] Details of Deployment and Presentation**

Deployment is in-process, as I have plans of creating a full-fledged website, which is an amazon-like respondent website.

## **[9] Conclusion and Future Scope**

The project was a new revelation but requires more fine touches to make it ready for heavy use. I plan on patenting the web scrapping algorithm and make it more inclusive to other smaller e-commerce platforms.

## **[10] References**

- <https://www.statista.com/statistics/379046/worldwide-retail-e-commerce-sales/>
- <https://huggingface.co/cardiffnlp/twitter-roberta-base-sentiment>
- <https://towardsdatascience.com/sentimental-analysis-using-vader-a3415fef7664>
- <https://www.geeksforgeeks.org/bag-of-words-bow-model-in-nlp/>
- <https://www.zenrows.com/blog/web-scraping-amazon>
- <https://www.datasciencecentral.com/how-to-scrape-amazon-product-data/>
- <https://www.youtube.com/watch?v=QpzMWQvxXWk>
- <https://www.youtube.com/watch?v=HiOtQMcl5wg>
- <https://medium.com/analytics-vidhya/web-scraping-amazon-reviews-a36bdb38b257>
- <https://www.geeksforgeeks.org/web-scraping-amazon-customer-reviews/>
- <https://www.digitalocean.com/community/tutorials/scrape-amazon-product-information-beautiful-soup>
- <https://www.geeksforgeeks.org/scraping-amazon-product-information-using-beautiful-soup/>
- <https://blog.apify.com/step-by-step-guide-to-scraping-amazon/>
- <https://github.com/AlexTheAnalyst/PortfolioProjects/blob/main/Amazon%20Web%20Scrape%20Project.ipynb>

## **Webscrapping.py Module**

```

1 # Importing required modules/libraries
2 from bs4 import BeautifulSoup
3 import requests
4 import pandas as pd
5 import numpy as np
6 import scipy.stats
7 import seaborn as sns
8 import random
9 import matplotlib.pyplot as plt
10 import warnings
11 import time
12 import datetime as dt
13 import plotly.express as px
14 import matplotlib.patheffects as pe
15 import regex as re
16 import smtplib
17 import re
18 from IPython.display import clear_output
19 warnings.filterwarnings("ignore")
20
21 # Finding user agent at: https://httpbin.org/get
22
23 # amazon orange, flipkart yellow
24 colors = ["#FEBD69", "#F8E831"]
25 customPalette = sns.set_palette(sns.color_palette(colors))
26
27 # Change some of seaborn's style settings with `sns.set()`
28 sns.set(style="ticks", # The 'ticks' style
29         rc={"figure.figsize": (6, 9), # width = 6, height = 9
30              "figure.facecolor": "ivory", # Figure colour
31              "axes.facecolor": "ivory"}, # Axes colour
32         palette=customPalette) # Axes colour
33
34 # dark black-blue amazon combo, flipkart blue
35 colors_2 = ["#37475A", "#047BD5"]
36
37 # Finding user agent at: https://httpbin.org/get
38
39
40 """
41 Amazon Web Scraping Class
42 """
43
44
45 class AmazonScraper:
46
47     def __init__(self, productName, pg_no):
48
49         self.user_query = productName # Storing the searched querry as an attribute
50
51         # Device Specific Headers
52         self.head = {
53             "User-Agent": "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/605.1.15 (KHTML, like
54             Gecko) "
55             "Version/16.1 Safari/605.1.15"}
56
57         # Driver Code as attributes (mandatory)
58         querry = self.generateUrl(productName, pg_no) # Getting the required querry
59         soup_scrapped = self.connectUrl(querry) # Connect to the amazon page and scrape it through
60         BeautifulSoup
61         # This will hold all the values of all
62         self.cells = soup_scrapped.find_all('div', {'class': 's-result-item',
63             'data-component-type': 's-search-result'})
64
65         # Creating a url request based on amazon's website url mechanisms
66         def generateUrl(self, query, pg_no):
67             query = re.sub(r"\s+", '+', query)
68             amazon_url = 'https://www.amazon.in/s?k={0}&page={1}'.format(query, pg_no)
69             return amazon_url
70
71         # connecting to the website
72         def connectUrl(self, url):
73             # s= requests.Session()
74             page = requests.get(url, headers=self.head)
75             soup = BeautifulSoup(page.content, "html.parser")
76             return soup
77
78         # Function to check whether no. of instances haven't exceeded 20
79         def lenCheck(self, arr):

```

```

78     return False
79
80 #         if len(arr)==20:
81 #             return
82 #         else:
83 #             return False
84
85 # Getter method to retrieve title(name) of each listed product
86 def getTitles(self):
87     product_titles = []
88
89     for results in self.cells:
90         item_name = results.h2.text # product title (static html --> h2)
91         product_titles.append(item_name)
92
93         if self.lenCheck(product_titles):
94             break
95
96     return product_titles
97
98 # Getter method to retrieve price(in rs.) of each listed product
99
100 def getPrices(self):
101     product_prices = []
102
103     for results in self.cells:
104         try:
105             item_price = results.find('span', {'class': "a-price-whole"}).text # (static html --> span
tag)
106             item_price = float(item_price.replace(",",""))
107
108         except AttributeError:
109             item_price = None
110
111         product_prices.append(item_price)
112
113         if self.lenCheck(product_prices):
114             break
115
116
117     return product_prices
118
119 # Getter method to retrieve review ratings of each listed product
120 def getRatings(self):
121     product_ratings = []
122
123     for results in self.cells:
124         try:
125             item_rating = results.find('i', {'class': 'a-icon'}).text
126
127             # Exception, if there are no reviews or ratings available
128             try:
129                 item_rating = float(item_rating[:3])
130
131             except ValueError:
132                 item_rating = 0
133
134             except AttributeError:
135                 item_rating = 0
136
137         product_ratings.append(item_rating)
138
139         if self.lenCheck(product_ratings):
140             break
141
142     return product_ratings
143
144 # Getter method to retrieve number of reviews of each listed product
145 def getRatingsCount(self):
146     product_ratings_count = []
147
148     for results in self.cells:
149         try:
150             item_rating_count = results.find('span', {'class': "a-size-base s-underline-text"}).text
151
152             # Exception, if there are no reviews or ratings available
153             try:
154                 item_rating_count = int(item_rating_count.replace(",",""))
155

```

```

156         except ValueError:
157             item_rating_count = 0
158
159     except AttributeError:
160         item_rating_count = 0
161
162     product_ratings_count.append(item_rating_count)
163
164     if self.lenCheck(product_ratings_count):
165         break
166
167     return product_ratings_count
168
169 # Getter method to retrieve images of each listed product
170 def getImages(self):
171     product_images = []
172     for results in self.cells:
173         try:
174             item_image = results.a.img["src"]
175
176         except TypeError:
177             item_image = None
178
179         except AttributeError:
180             item_image = None
181
182         product_images.append(item_image)
183
184     if self.lenCheck(product_images):
185         break
186
187     return product_images
188
189 # Getter method to retrieve urls of each listed product
190 def getUrl(self):
191     product_urls = []
192
193     for results in self.cells:
194         try:
195             item_url = "https://www.amazon.in" + results.h2.a["href"]
196
197         except AttributeError:
198             item_url = None
199
200         product_urls.append(item_url)
201
202     if self.lenCheck(product_urls):
203         break
204
205     return product_urls
206
207 # Converting the scrapped data into an appropriate dataframe
208 def toDataframe(self):
209
210     # Printing the array lengths while testing for discrepancies
211     """
212     print(len(self.getTitles()))
213     print(len(self.getUrl()))
214     print(len(self.getPrices()))
215     print(len(self.getRatings()))
216     print(len(self.getRatingsCount()))
217     print(len(self.getImages()))
218     """
219
220     df = pd.DataFrame({
221         "query_searched": self.user_query,
222         "ecommerce_website": "Amazon",
223         "product_name": self.getTitles(),
224         "product_price": self.getPrices(),
225         "product_rating": self.getRatings(),
226         "rating_count": self.getRatingsCount(),
227         "product_image_url": self.getImages(),
228         "product_url": self.getUrl()
229     })
230
231     return df
232
233     # Additional function to go to each product's url and extract brand name, request response time too
234     long

```

```

234     def addBrand(self, df):
235         product_brands = []
236
237         for x in range(len(df)):
238             url = df.loc[x][7]
239             page = requests.get(url, headers=self.head)
240             soup = BeautifulSoup(page.content, "html.parser")
241
242             try:
243                 product_overview = soup.find(id="productOverview_feature_div")
244                 cell_detailed = product_overview.find('div', {"class": "a-section a-spacing-small a-spacing-top-small"})
245                 cells = cell_detailed.find_all('span')
246
247                 y = "Brand"
248                 for x in range(0, len(cells)):
249                     if y in cells[x]:
250                         item_brand = cells[x + 1].text
251
252             except AttributeError:
253                 item_brand = 0
254
255             product_brands.append(item_brand)
256
257         df["product_brand"] = products_brands
258
259
260
261 """
262 Flipkart Web Scraping Class
263 """
264
265
266 class FlipkartScraper:
267
268     def __init__(self, productName, pg_no):
269
270         self.grid = False
271         self.img = True
272
273         self.user_query = productName # Storing the searched query as an attribute
274
275         # Device Specific Headers
276         self.head = {
277             "User-Agent": "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/605.1.15 (KHTML,
278             like Gecko) "
279             "Version/16.1 Safari/605.1.15"
280
281             # Driver Code as attributes (mandatory)
282             querry = self.generateUrl(productName, pg_no) # Getting the required query
283             self.soup_scrapped = self.connectUrl(querry) # Connect to the Flipkart page and scrape it through
284             BeautifulSoup
285
286             # This will hold all the values of all cells in the page
287             self.cells = self.soup_scrapped.find_all('a', {"class": "_1fQZEK"})
288
289             # Flagging self.grid incase the layout of flipkart webpage is gridwise, instead of row wise
290             if len(self.cells) == 0:
291                 self.grid = True
292                 self.cells = self.soup_scrapped.find_all('div', {"class": "_4ddWXP"})
293
294             if len(self.cells) == 0:
295                 self.grid = True
296                 self.img = False
297                 self.cells = self.soup_scrapped.find_all('div', {"class": "_1xH6tK _373qXS"})
298
299             else:
300                 self.grid = False
301
302             # Creating a url request based on amazon's website url mechanisms
303             def generateUrl(self, query, pg_no):
304                 query = re.sub(r"\s+", '+', query)
305                 flipkart_url = 'https://www.flipkart.com/search?q={0}&page={1}'.format(query, pg_no)
306                 return flipkart_url
307
308             # connecting to the website
309             def connectUrl(self, url):
310                 page = requests.get(url, headers=self.head)
311                 soup = BeautifulSoup(page.content, "html.parser")

```

```

310     return soup
311
312     # Getter method to retrieve title(name) of each listed product
313     def getTitles(self):
314         product_titles = []
315
316         for index in range(len(self.cells)):
317             try:
318                 if self.grid:
319                     item_title = self.cells[index].find_all('a')[1]['title']
320
321                 else:
322                     item_title = self.cells[index].find('div', {"class": "_4rR01T"}).text
323
324             except AttributeError:
325                 item_title = ''
326
327             except KeyError:
328                 item_title = ''
329
330             except ValueError:
331                 item_title = ''
332
333             product_titles.append(item_title)
334
335     return product_titles
336
337     # Getter method to retrieve price(in rs.) of each listed product
338     def getPrices(self):
339         product_prices = []
340         for index in range(len(self.cells)):
341             try:
342                 if self.grid:
343                     str_price = self.cells[index].find_all('a')[2].find('div', {"class": "_30jeq3"}).text
344                     item_price = float(str_price[1:].replace(","))
345
346                 else:
347                     item_price = (self.cells[index].find('div', {"class": "_30jeq3 _1_WHN1"}).text)
348                     item_price = float((item_price[1:]).replace(","))
349
350             except AttributeError:
351                 item_price = None
352
353             product_prices.append(item_price)
354
355     return product_prices
356
357     # Getter method to retrieve review ratings of each listed product
358     def getRatings(self):
359         product_ratings = []
360         for index in range(len(self.cells)):
361
362             try:
363                 if self.grid:
364                     item_rating = float(self.cells[index].find('div', {"class": '_3LWZLK'}).text)
365
366                 else:
367                     item_rating = float(self.cells[index].find('div', {"class": "_3LWZLK"}).text)
368
369             except AttributeError:
370                 item_rating = 0
371
372             except ValueError:
373                 item_rating = 0
374
375             product_ratings.append(item_rating)
376
377     return product_ratings
378
379     # Getter method to retrieve number of ratings of each listed product
380     def getRatingsCount(self):
381         product_ratings_count = []
382         for index in range(len(self.cells)):
383
384             try:
385
386                 if self.grid:
387                     str_ratings_count = self.cells[index].find('span', {'class': '_2_R_DZ'}).text
388                     item_rating_count = int(str_ratings_count.replace("(", "").replace(")", ""))

```



```

466             except AttributeError:
467                 item_url = None
468
469             except KeyError:
470                 item_url = None
471
472         product_urls.append(item_url)
473
474     return product_urls
475
476
477 # Converting the scrapped data into an appropriate dataframe
478 def toDataframe(self):
479
480     # Printing the array lengths while testing for discrepancies
481     """
482     print(len(self.getTitles()))
483     print(len(self.getUrls()))
484     print(len(self.getPrices()))
485     print(len(self.getRatings()))
486     print(len(self.getRatingsCount()))
487     print(len(self.getReviewsCount()))
488     print(len(self.getImages()))
489     """
490
491     df = pd.DataFrame({
492         "query_searched": self.user_query,
493         "ecommerce_website": "Flipkart",
494         "product_name": self.getTitles(),
495         "product_price": self.getPrices(),
496         "product_rating": self.getRatings(),
497         "rating_count": self.getRatingsCount(),
498         "product_image_url": self.getImages(),
499         "product_url": self.getUrls(),
500         "reviews_count": self.getReviewsCount()
501     })
502
503     return df
504
505
506 """
507 Comments Scraping (Extraction) (Currently Amazon Specific)
508 """
509
510
511 class ScrapeComments:
512
513     def __init__(self, df, product_qrd, max_reviews):
514
515         # Device Specific Headers
516         self.head = {
517             "User-Agent": "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/605.1.15 (KHTML, like Gecko) Version/16.1 Safari/605.1.15"
518
519             "Version/16.1 Safari/605.1.15"
520
521         self.df_comments = self.getAmazonComments(df, product_qrd, max_reviews)
522
523     # Function that returns the Data Frame created (extracted comments)
524     def getDF(self):
525         return self.df_comments
526
527     # Main Function that extracts and converts the comments into a structured dataframe
528     def getAmazonComments(self, df, onclick_name, max_reviews):
529
530         # Counters
531         page_no = 0 # All reviews --> current page number counter
532         reviews = 0 # Total reviews extracted counter
533         reviews_count = 0 # Counter initialised to prevent item called before assignment error
534
535         # Empty lists that will hold all reviewee data like name, date of review, ratings given, review
536         title, content
537         customer_name = []
538         review_date = []
539         ratings = []
540         review_title = []
541         review_content = []
542         # global reviews_soup
543
544         # Finding the url to the queried product name using indexing

```

```

543     item_url = list(df['product_url'][df['product_name'] == onclick_name])[0]
544
545     # Connecting (sending a request) to the product's page
546     item_page = requests.get(item_url, headers=self.head)
547     item_soup = BeautifulSoup(item_page.content, "html.parser")
548
549     try:
550
551         # Till the reviews counter is less than the reviews_count mentioned on the 'all reviews' page
552         while reviews <= reviews_count:
553
554             # Finding the link to "all reviews on the product's page", and adding iterating page_no
555             reviews_url = item_soup.find('a', {'class': 'a-link-emphasis a-text-bold'})['href']
556             reviews_url = "https://www.amazon.in" + reviews_url + "&pageNumber=" + str(page_no)
557
558             # Sending request for the Reviews Url page
559             reviews_page = requests.get(reviews_url, headers=self.head)
560             reviews_soup = BeautifulSoup(reviews_page.content, "html.parser")
561
562             # Number of reviews for the particular page, same throughout all comment (all reviews)
563             reviews_count = reviews_soup.find('div', {'class': "a-row a-spacing-base a-size-base"}).
564             text
565             reviews_count = int(
566                 (reviews_count.strip().split('ratings,'))[1].split('with')[0].strip().replace(",",""))
567
568             # Extracting the required comments in a page
569             # These are generally lists containing 10 review related data each.
570             names = reviews_soup.select('span.a-profile-name')[2:]
571             titles = reviews_soup.select('a.review-title span')
572             dates = reviews_soup.select('span.review-date')[2:]
573             stars = reviews_soup.select('i.review-rating span.a-icon-alt')[2:]
574             content = reviews_soup.select('span.review-text-content span')
575
576             # IndexError Handling : If the sections are left empty by the customers, the index is out
577             # of range
578             for count in range(len(dates)):
579
580                 try:
581                     customer_name.append(names[count].get_text())
582                 except IndexError:
583                     customer_name.append(0)
584
585                 try:
586                     review_date.append(dates[count].get_text().replace("Reviewed in India on ", ""))
587                 except IndexError:
588                     review_date.append(0)
589
590                 try:
591                     review_title.append(titles[count].get_text())
592                 except IndexError:
593                     review_title.append(None)
594
595                 try:
596                     ratings.append(float(stars[count].get_text()[:2]))
597                 except IndexError:
598                     ratings.append(None)
599
600                 try:
601                     review_content.append(content[count].get_text())
602                 except IndexError:
603                     review_content.append(None)
604
605             reviews = reviews + 1
606
607             # Incase the number of reviews counter has reached the maximum reviews to be extracted
608             if reviews == max_reviews:
609                 break
610
611             page_no = page_no + 1
612             # print(page_no)
613
614         except:
615             pass
616
617         """
618         Debugging Section

```

```

618     print(len(customer_name))
619     print(len(review_date))
620     print(len(review_title))
621     print(len(ratings))
622     print(len(review_content))
623     """
624
625     # Recalling required functions from the module.
626     df = self.commentsToDf(customer_name, review_date, review_title, ratings, review_content)
627     df = self.cleanReviewDates(df)
628
629     return df
630
631     # Function that assigns the lists extracted to a dataframe and returns it.
632     def commentsToDf(self, names, dates, titles, ratings, content):
633
634         reviews_df = pd.DataFrame({ "customer_name": names,
635                                     "review_date": dates,
636                                     "review_title": titles,
637                                     "review_ratings": ratings,
638                                     "review_content": content})
639
640         return reviews_df
641
642     # Function to convert dates which have indian flag and incorrect redundant formats in it, st object in
643     # dd/mm/yy
644     def cleanReviewDates(self, df):
645
646         cleaned_dates = []
647         for comment in list(df['review_date']):
648             for string in range(len(comment)):
649                 if comment[string].isnumeric():
650                     cleaned_dates.append(comment[string:])
651                     break
652
653         df['review_date'] = cleaned_dates
654         df['review_date'] = pd.to_datetime(df['review_date']).dt.strftime('%d/%m/%Y')
655
656         return df
657
658 """
659 Functions for Creating and merging datasets acccording to a user's querry
660 """
661
662
663 # Function that concats 'n' dataframes
664 def joinDf(df_arr):
665     df = pd.DataFrame()
666     for x in df_arr:
667         df = pd.concat([df, x],
668                       ignore_index=True)
669
670     return df
671
672 # User Defined Funciton : that checks whether the brand is present in the product_name
673 def brandCheck(df, brand):
674     if len(brand) == 0:
675         return df
676
677     else:
678         # Creating a regex expression so that if there's upper or lower case in brand or product it'll be
679         accounted
680         letters = "("
681         for letter in brand:
682             upper_x = letter.upper()
683             lower_x = letter.lower()
684             letters = letters + "[" + upper_x + lower_x + "]"
685
686         letters = letters + ")*\s"
687         # print(letters)
688
689         arr_vals = []
690         for x in range(0, len(df)):
691             product_name = str(df.loc[x][2])
692             index_brand = re.search(letters, product_name)
693             if index_brand:
694                 arr_vals.append(df.loc[x])

```

```

695         df = pd.DataFrame(arr_vals, columns=list(df.columns))
696
697     return df
698
699
700 # Function that combines the amazon and filpkart webscrapped
701 def finalDf(amz_df, flk_df, brand_name):
702     arr_dfs = []
703     arr_dfs.append(amz_df)
704     arr_dfs.append(flk_df)
705     df1 = joinDf(arr_dfs)
706     df = brandCheck(df1, brand_name)
707     # df = removePercentile(df,20)
708     df.reset_index(inplace=True, drop=True)
709     # df.drop(['index'], axis=1, inplace=True)
710
711     return df1, df
712
713
714 # MAIN FUNCTION that takes in the usr query
715 def execCode(query, brand):
716     # Initila page number
717     pg_no = 1
718     amz_df_arr = []
719     flk_df_arr = []
720
721     # A variable that creates a 100 second buffer between current time and the end.
722     t_end = time.time() + 100
723
724     # Until the current sys time isn't less than the end time variable
725     while time.time() < t_end:
726         try:
727             amz_1 = AmazonScraper(query, pg_no)
728             flk_1 = FlipkartScraper(query, pg_no)
729
730             amz_df = amz_1.toDataframe()
731             flk_df = flk_1.toDataframe()
732
733             amz_df_arr.append(amz_df)
734             flk_df_arr.append(flk_df)
735
736             pg_no = pg_no + 1
737
738
739         except:
740             break
741
742
743
744     # Joining both the amazon and flipkart pages individually
745     amz_df = joinDf(amz_df_arr)
746     flk_df = joinDf(flk_df_arr)
747
748     df, brand_present_df = finalDf(amz_df, flk_df, brand)
749
750     return df, brand_present_df
751

```

## **ProjectVisualisations.py Module**

```

1 # Importing required libraries for further analysis
2 import pandas as pd
3 import numpy as np
4 import scipy.stats
5 import seaborn as sns
6 import random
7 import matplotlib.pyplot as plt
8 import warnings
9 import time
10 import plotly.express as px
11 import matplotlib.patheffects as pe
12 import regex as re
13
14 warnings.filterwarnings("ignore")
15
16 # Finding user agent at: https://httpbin.org/get
17
18 # amazon orange, flipkart yellow
19 colors = ["#FEBD69", "#F8E831"]
20 customPalette = sns.set_palette(sns.color_palette(colors))
21
22 # Change some of seaborn's style settings with `sns.set()`
23 sns.set(style="ticks", # The 'ticks' style
24         rc={"figure.figsize": (6, 9), # width = 6, height = 9
25              "figure.facecolor": "ivory", # Figure colour
26              "axes.facecolor": "ivory"}, # Axes colour
27         palette=customPalette) # Axes colour
28
29 # dark black-blue amazon combo, flipkart blue
30 colors_2 = ["#37475A", "#047BD5"]
31
32 from webscrapping import AmazonScraper
33 from webscrapping import FlipkartScraper
34 from webscrapping import ScrapeComments
35
36
37 # Function to count products with brand name present in comparison to the original web - scrapped dataset
38 def countActualProductsPlot(df, df_correct):
39     # if len(df)==len(df_correct):
40     #     myexplode=[0.5,0]
41     #     plt.pie(x = [100,0], explode = myexplode, shadow = True,
42     #             colors = ["g", "r"])
43     #     plt.legend(["Relevant Products", "Irrelevant Products"])
44     #     plt.title("No brand mentioned / All Products Have Brand Present in Name",
45     #               fontsize=30,color="b")
46     #     return plt.show()
47     # else:
48     try:
49         amazon_prod_count = sum(df["ecommerce_website"] == "Amazon")
50         flipkart_prod_count = sum(df["ecommerce_website"] == "Flipkart")
51
52         actual_amazon_prod_count = sum(df_correct["ecommerce_website"] == "Amazon")
53         actual_flipkart_prod_count = sum(df_correct["ecommerce_website"] == "Flipkart")
54
55         print("-> Amazon has",
56               amazon_prod_count,
57               "products listed on each page for the queried product,out of which ",
58               (amazon_prod_count - actual_amazon_prod_count),
59               " are of other brands")
56
56         print("-> Flipkart has",
57               flipkart_prod_count,
58               "products listed on each page for the queried product,out of which ",
59               actual_flipkart_prod_count,
60               " are of other brands")
66
67         plt.figure(figsize=(12, 7))
68         plt.rcParams["axes.edgecolor"] = "black"
69         plt.rcParams["axes.linewidth"] = 2.50
70         plt.suptitle("Stacked Bar Chart Comparing Relevance of Listed Products", fontsize=30, color='b')
71
72         # Bar Plot
73         plt.subplot(1, 2, 1)
74         x = ["Amazon", "Flipkart"]
75         y1 = [amazon_prod_count - actual_amazon_prod_count,
76               flipkart_prod_count - actual_flipkart_prod_count]
77
78         y2 = [actual_amazon_prod_count, actual_flipkart_prod_count]
79

```

```

80
81     # plot bars in stack manner
82     plt.bar(x, y1, color='red', edgecolor=colors_2)
83     plt.bar(x, y2, bottom=y1, color=colors, edgecolor=colors_2)
84
85     # Labels
86     plt.xlabel("E-Commerce Website", fontsize=16, color="red")
87     plt.ylabel("Count of Products", fontsize=16, color="red")
88     plt.legend(["Irrelevant Products", "Relevant Products Amazon", "Relevant Products Flipkart"])
89
90     # Pie Chart
91     plt.subplot(1, 2, 2)
92     myexplode = [0.4, 0]
93
94     relv_prd = sum(y2) / (amazon_prod_count + flipkart_prod_count)
95     irrelv_prd = 1 - relv_prd
96
97     plt.pie(x=[relv_prd, irrelv_prd], explode=myexplode, shadow=True,
98             colors=["g", "r"])
99     plt.legend(["Relevant Products", "Irrelevant Products"])
100    plt.tight_layout()
101    return plt.show()
102
103 except:
104     print("Error !")
105
106
107 def countProductsPlot(df):
108     amazon_prod_count = sum(df["ecommerce_website"]=="Amazon")
109     flipkart_prod_count = sum(df["ecommerce_website"] == "Flipkart")
110
111     print("-> Amazon has", amazon_prod_count, "products listed on each page.")
112     print("-> Flipkart has", flipkart_prod_count, "products listed on each page.")
113
114     plt.figure(figsize=(12, 7))
115     plt.rcParams["axes.edgecolor"] = "black"
116     plt.rcParams["axes.linewidth"] = 2.50
117     plt.suptitle("Products Listed Per Page for the Query", fontsize=30, color='b')
118     plt.subplot(1, 2, 1)
119     # fig.subtitle('Products per page for given product')
120
121     sns.countplot(x='ecommerce_website', data=df, color=customPalette,
122                   linewidth=3, edgecolor=colors_2)
123
124     plt.xlabel("Ecommerce Website", fontsize=16, color="red")
125     plt.ylabel("No. of Products per page", fontsize=16, color="red")
126
127     plt.subplot(1, 2, 2)
128     plt.pie(x=[amazon_prod_count,flipkart_prod_count], labels=["Amazon", "Flipkart"],
129             colors=customPalette, autopct='%.0f%')
130
131     plt.tight_layout()
132     return plt.show()
133
134
135 # Function to find average of a column for the particular e-commerce website
136 def subset_avg_gen(df, website, col):
137     df_subset = df[df["ecommerce_website"] == website]
138     average = df_subset[col].mean()
139     return average
140
141
142 # User Defined Function : Average Price of the products displayed upon query, website wise
143 def priceCompLine(df):
144     x = list(df[df["ecommerce_website"] == "Amazon"]["product_price"])
145     y = list(df[df["ecommerce_website"] == "Flipkart"]["product_price"])
146     plt.figure(figsize=(15, 10))
147     plt.rcParams["axes.edgecolor"] = "black"
148     plt.rcParams["axes.linewidth"] = 2.50
149     overlapping = 0.1
150
151     plt.plot(x, lw=6, path_effects=[pe.SimpleLineShadow(shadow_color=colors_2[0]), pe.Normal()])
152     plt.plot(y, lw=6, path_effects=[pe.SimpleLineShadow(shadow_color=colors_2[1]), pe.Normal(), ])
153
154     plt.legend(["Amazon", "Flipkart"], loc='upper left')
155     plt.xlabel("Product No.", color="red")
156     plt.ylabel("Price (in rs.)", color="red")
157     plt.title("Comparing Prices (line-plot)", fontsize=30, c='b')
158     return plt.show()

```

```

159
160
161 def CompPriceBox(df):
162     amz_avg_price = subset_avg_gen(df, "Amazon", "product_price")
163     flk_avg_price = subset_avg_gen(df, "Flipkart", "product_price")
164
165     print("Average Product Price for queried product at amazon is : ", amz_avg_price)
166     print("Average Product Price for queried product at flipkart is : ", flk_avg_price)
167
168     plt.rcParams["axes.edgecolor"] = "black"
169     plt.rcParams["axes.linewidth"] = 2.50
170     plt.figure(figsize=(15, 10))
171     # Box plot
172     b = sns.boxplot(data=df,
173                      x="ecommerce_website",
174                      y="product_price",
175                      width=0.4,
176                      linewidth=2,
177
178                      showfliers=False)
179     # Strip plot
180
181     b = sns.stripplot(data=df,
182                        x="ecommerce_website",
183                        y="product_price",
184                        # Colours the dots
185                        linewidth=1,
186                        alpha=0.4)
187
188     b.set_ylabel("Price", fontsize=14, color="red")
189     b.set_xlabel("E-Commerce Website", fontsize=14, color="red")
190     b.set_title("Comparing Prices (box-plot)", fontsize=30, c='b')
191
192     sns.despine(offset=5, trim=True)
193
194     return b.get_figure()
195
196
197 # User Defined Function : Number of products available on Amazon and flipkart for the given querry
198 def priceNamePlotly(df):
199     fig = px.line(df, x="product_name", y="product_price",
200                   color="ecommerce_website", template="plotly_dark",
201                   color_discrete_sequence=colors,
202                   title="Products and their price")
203
204     fig.update_xaxes(showticklabels=False)
205     return fig.show()
206
207
208 # User Defined Function : Average Rating of the products displayed upon query, website wise
209 def CompRatingBox(df):
210     amz_avg_ratings = subset_avg_gen(df, "Amazon", "product_rating")
211     flk_avg_ratings = subset_avg_gen(df, "Flipkart", "product_rating")
212
213     print("Average Product Rating for queried product at amazon is : ", amz_avg_ratings)
214     print("Average Product Rating for queried product at flipkart is : ", flk_avg_ratings)
215
216     plt.figure(figsize=(15, 10))
217
218     # Box plot
219     b = sns.boxplot(data=df,
220                      x="ecommerce_website",
221                      y="product_rating",
222                      width=0.4,
223                      linewidth=2,
224
225                      showfliers=False)
226     # Strip plot
227     b = sns.stripplot(data=df,
228                        x="ecommerce_website",
229                        y="product_rating",
230                        # Colours the dots
231                        linewidth=1,
232                        alpha=0.4)
233
234     b.set_ylabel("Ratings (out of 5)", fontsize=14, color="red")
235     b.set_xlabel("E-Commerce Website", fontsize=14, color="red")
236     b.set_title("Comparing Ratings (box-plot)", fontsize=30, c='b')
237

```

```

238     sns.despine(offset=5, trim=True)
239     # b.get_figure()
240     return b
241
242
243 # User Defined Function : Analysing the relationship between Product's Price and its Rating
244 # Function that plots a scatter plot between price and ratings
245 def price_rat_corr(df):
246     # Printing what kind of correlation is present between product price and product rating
247     price_rating_corr = df.corr()['product_rating']['product_price']
248
249     if abs(price_rating_corr) > 0.9:
250         print("Product Ratings and Product Prices are very highly correlated. ")
251
252     if abs(price_rating_corr) > 0.7:
253         print("Product Ratings and Product Prices are highly correlated. ")
254
255     if abs(price_rating_corr) > 0.5:
256         print("Product Ratings and Product Prices are moderately correlated. ")
257
258     if abs(price_rating_corr) > 0.3:
259         print("Product Ratings and Product Prices have low correlation. ")
260
261 else:
262     print("Product Ratings and Product Prices have very little to no correlation. ")
263
264 fig = px.scatter(df, x="product_rating",
265                   color="ecommerce_website", template="plotly_dark",
266                   color_discrete_sequence=colors,
267                   trendline="ols",
268                   title="Analysing the relationship between Product's Price and Product's Ratings")
269
270 fig.update_xaxes(showticklabels=True)
271 return fig.show()
272
273
274 """
275 Main function for visualising all the results
276 """
277
278
279 def visualiseQueryResults(df, brand_present_df):
280     # amazon orange, flipkart yellow
281     colors = ["#FEBD69", "#F8E831"]
282     customPalette = sns.set_palette(sns.color_palette(colors))
283
284     # Change some of seaborn's style settings with `sns.set()`
285     sns.set(style="ticks", # The 'ticks' style
286             rc={ "figure.facecolor": "ivory", # Figure colour
287                  "axes.facecolor": "ivory"}, # Axes colour
288             palette=customPalette) # Axes colour
289
290     ax1 = countProductsPlot(df)
291     if len(brand_present_df)==len(df):
292         pass
293     else:
294         ax2 = countActualProductsPlot(df, brand_present_df)
295
296     ax3 = priceComplLine(brand_present_df)
297     ax4 = CompPriceBox(brand_present_df)
298     ax5 = priceNamePlotly(brand_present_df)
299
300     ax6 = CompRatingBox(brand_present_df)
301
302     ax7 = price_rat_corr(brand_present_df)
303

```

## **ProductReviewAnalysis.py Module**

```

1 # Importing Required libraries and modules
2
3 from bs4 import BeautifulSoup
4 import requests
5 import time
6 import datetime
7
8 import smtplib
9 import re
10
11 import pandas as pd
12 import seaborn as sns
13 import datetime as dt
14 import numpy as np
15 # Finding user agent at: https://httpbin.org/get
16
17 from nltk.sentiment import SentimentIntensityAnalyzer
18 from tqdm.notebook import tqdm
19
20 from nltk.sentiment import SentimentIntensityAnalyzer
21 from tqdm.notebook import tqdm
22 from transformers import AutoTokenizer
23 from transformers import AutoModelForSequenceClassification
24 from scipy.special import softmax
25 from IPython.display import clear_output
26 from tabulate import tabulate
27 import urllib.request
28 # from PIL import Image
29 from IPython.display import Image
30 from wordcloud import WordCloud, STOPWORDS
31 import matplotlib.pyplot as plt
32
33 stopwords = set(STOPWORDS)
34
35 # Change some of seaborn's style settings with `sns.set()`
36 sns.set(style="ticks", # The 'ticks' style
37          rc={"figure.figsize": (6, 9), # width = 6, height = 9
38               "figure.facecolor": "ivory", # Figure colour
39               "axes.facecolor": "ivory"}) # Axes colour
40
41 # Importing files created by me, and saved in. webscraping.py
42 from webscraping import AmazonScraper
43 from webscraping import FlipkartScraper
44 from webscraping import ScrapeComments
45
46 MODEL = "cardiffnlp/twitter-roberta-base-sentiment"
47 tokenizer = AutoTokenizer.from_pretrained(MODEL)
48 model = AutoModelForSequenceClassification.from_pretrained(MODEL)
49
50 class color:
51     PURPLE = '\033[95m'
52     CYAN = '\033[96m'
53     DARKCYAN = '\033[36m'
54     BLUE = '\033[94m'
55     GREEN = '\033[92m'
56     YELLOW = '\033[93m'
57     RED = '\033[91m'
58     BOLD = '\033[1m'
59     UNDERLINE = '\033[4m'
60     END = '\033[0m'
61
62
63 """
64 User Defined Functions
65 """
66
67
68 # User Defined Function To neatly print the products of the particular dataset
69 def printTableProducts(df, length):
70     products_dict = dict(df[df["ecommerce_website"] == "Amazon"]["product_name"])
71     arr = []
72
73     for x in products_dict:
74         arr.append([x, products_dict.get(x)[:70]])
75
76     if len(arr) == length:
77         break
78
79     print(tabulate(arr, headers=['Product No.', 'Product Name']))

```

```

80
81
82 # User Defined Function To Show the selected products details in a neat manner
83 def showProductDeets(df, index):
84     element = df.iloc[index]
85
86     print((color.BOLD + color.BLUE + 'PRODUCT DETAILS' + color.END).center(100, "-"))
87     print()
88     print(color.BOLD + 'Product Name : ' + color.END, element[2])
89     print(color.BOLD + 'Product Price : ' + color.END, element[3])
90     print(color.BOLD + 'Average Product Rating : ' + color.END, element[4])
91     print(color.BOLD + 'Total Number of Ratings : ' + color.END, element[5])
92     print()
93
94     urllib.request.urlretrieve(df.iloc[index][6], "ex.png")
95     display(Image(filename='ex.png'))
96     print("-" * 100)
97
98
99 # Function to add a new column that contains the size of each review
100 def addReviewLen(df):
101     len_review = []
102     for x in range(len(df)):
103         rev_size = len(df.iloc[x][4])
104         len_review.append(rev_size)
105     df['review_size'] = len_review
106
107
108 # Function to clean the given comments dataset (drop na/null values)
109 def cleanCommentsDf(df):
110     df.dropna(inplace=True)
111     df.reset_index(inplace=True, drop=True)
112
113
114 # Return Mapped value based on Categorical variables for weighted calculations
115 def mapValue(element):
116     # 'Positive', 'Neutral', 'Negative' unique values
117     if element == 'Positive':
118         return 1
119
120     if element == 'Neutral':
121         return 0
122
123     else:
124         return -1
125
126
127 """
128 VADER AND ROBERTA MODEL related user defined functions
129 """
130
131
132 # Vader model (applying vader model on each comment and creating a df)
133 def getVaderDf(df):
134     global sia
135     sia = SentimentIntensityAnalyzer()
136     # Runn the polarity score on the entire dataset
137     res = {}
138     for i, row in tqdm(df.iterrows(), total=len(df), desc="Vader Model"):
139         text = row['review_content']
140         index = list(df.index)[i]
141         val = sia.polarity_scores(text)
142         res[index] = {
143             'vader_neg': float(val['neg']),
144             'vader_pos': float(val['pos']),
145             'vader_comp': float(val['compound'])
146         }
147
148     # Converting thed dictionary into a dataframe and using transpose to make columns as rows
149     vaders = pd.DataFrame(res).T
150     vaders.reset_index(inplace=True, drop=True)
151     # Merfing the vaders dataframe with the original dataset
152
153     return vaders
154
155
156 # Roberta Model Scores (individual)
157 def polarity_scores_roberta(example):
158     if len(example)>2150:

```

```

159     example = example[:2000]
160     #if len(example.split()) > 500:
161     #example = " ".join(example.split()[:500])
162
163     encoded_text = tokenizer(example, return_tensors='pt')
164     output = model(**encoded_text)
165     scores = output[0][0].detach().numpy()
166     scores = softmax(scores)
167     scores_dict = {
168         'roberta_neg': scores[0],
169         'roberta_neu': scores[1],
170         'roberta_pos': scores[2]
171     }
172     return scores_dict
173
174
175 # Roberta model (applying Roberta model on each comment and creating a df)
176 def getRobertaDf(df):
177
178     # Runn the polarity score on the entire dataset
179     res = {}
180     for i, row in tqdm(df.iterrows(), total=len(df), desc="Roberta Model"):
181         text = row['review_content']
182         index = list(df.index)[i]
183         val = polarity_scores_roberta(text)
184         res[index] = val
185
186     # Converting thed dictionary into a dataframe and using transpose to make columns as rows
187     df = pd.DataFrame(res).T
188     df.reset_index(inplace=True, drop=True)
189
190     return df
191
192
193 # User defined function to create a categorical column on the basis of the three Vader values
194 def getVaderSentiment(df):
195     arr = []
196     for x in range(len(df)):
197         vader_compound_score = df["vader_comp"][x]
198
199         if vader_compound_score >= 0.05:
200             arr.append("Positive")
201
202         elif vader_compound_score <= - 0.05:
203             arr.append("Negative")
204
205         else:
206             arr.append("Neutral")
207
208     if 'vader_sentiment' not in df.columns:
209         df.insert(loc=6,
210                 column='vader_sentiment',
211                 value=arr)
212
213     else:
214         df['vader_sentiment'] = arr
215
216
217 # User defined function to create a categorical column on the basis of the three Roberta models
218 def getRobertaSentiment(df):
219     arr = []
220     for pos in range(len(df)):
221         roberta_sentiment_dict = {"Positive": df[["roberta_pos", "roberta_neg", "roberta_neu"]].iloc[pos][0],
222 },                                     "Negative": df[["roberta_pos", "roberta_neg", "roberta_neu"]].iloc[pos][1],
223 ],                                     "Neutral": df[["roberta_pos", "roberta_neg", "roberta_neu"]].iloc[pos][2]
224
225         roberta_sentiment_value = max(roberta_sentiment_dict, key=roberta_sentiment_dict.get)
226
227         arr.append(roberta_sentiment_value)
228
229     if 'roberta_sentiment' not in df.columns:
230         df.insert(loc=6,
231                 column='roberta_sentiment',
232                 value=arr)
233
234     else:

```

```

235         df['roberta_sentiment'] = arr
236
237
238 # User defined function : Creating a weighted column based on
239 def getFinalSentiment(df):
240     arr = []
241     roberta_sentiment_vals = list(map(mapValue, list(df['roberta_sentiment'])))
242     vader_sentiment_vals = list(map(mapValue, list(df['vader_sentiment'])))
243
244     roberta_val = list(df['roberta_sentiment'])
245     vader_val = list(df['vader_sentiment'])
246
247     for x in range(len(df)):
248         val = float((roberta_sentiment_vals[x] * 0.6) + (vader_sentiment_vals[x] * 0.4))
249
250         if val >= 0.19:
251             arr.append("Positive")
252
253         elif val >= (-0.4):
254             arr.append("Neutral")
255
256
257         elif val >= (-1):
258             arr.append("Negative")
259
260     if 'sentiment' not in df.columns:
261         df.insert(loc=6,
262                   column='sentiment',
263                   value=arr)
264
265     else:
266         df['sentiment'] = arr
267
268
269 """
270 Plotting related user defined functions
271 """
272
273
274 # Function for plotting a comparison on the kinds of reviews given for the product
275 def ratingsItemPlot(df):
276     plt.figure(figsize=(12, 7))
277     plt.subplot(1, 2, 1)
278
279     rev_count = list(df['review_ratings'].value_counts().sort_index())
280     rev_tags = list(df['review_ratings'].unique())
281     rev_tags.sort()
282
283     plt.bar(rev_tags, rev_count)
284
285     plt.plot(rev_tags, rev_count, color='black', linewidth=7.0)
286
287     plt.xlabel('Review Stars', fontsize=16, color="red")
288     plt.ylabel('Review Count', fontsize=16, color="red")
289
290     plt.suptitle("Reviews Count for the product", fontsize=30, color='b')
291
292     plt.subplot(1, 2, 2)
293     #myexplode = [0.15, 0.19, 0.22, 0.12, 0.15]
294
295     plt.pie(x=rev_count,
296             shadow=True, autopct='%.0f%')
297
298     plt.legend(rev_tags)
299     # plt.rcParams.update({'font.size': 10})
300
301     plt.tight_layout()
302     return plt.show()
303
304
305 # User defined function to count and show a woordle like frequency chart
306 def show_wordcloud(data):
307     wordcloud = WordCloud(
308         background_color='white',
309         stopwords=stopwords,
310         max_words=200,
311         max_font_size=40,
312         scale=3,
313         random_state=1 # chosen at random by flipping a coin; it was heads

```

```

314     ).generate(str(data))
315
316     fig = plt.figure(1, figsize=(15, 15))
317     plt.axis('off')
318
319     plt.imshow(wordcloud)
320
321     return plt.show()
322
323 # Graph that compares the review length and the ratings given
324 def sizeComPlot(df):
325     plt.figure(figsize=(12, 7))
326     plt.suptitle("Product Rating vs Review Length", fontsize=40, color="blue")
327     plt.subplot(1, 2, 1)
328     sns.pointplot(data=df, x="review_ratings", y="review_size")
329     plt.xlabel("Rating", fontsize=16, color="red")
330     plt.ylabel("Review Length", fontsize=16, color="red")
331
332     plt.subplot(1, 2, 2)
333     sns.boxplot(data=df, x="review_ratings", y="review_size")
334     plt.xlabel("Rating", fontsize=16, color="red")
335     plt.ylabel("Review Length", fontsize=16, color="red")
336     plt.tight_layout()
337
338     return plt.show()
339
340 # User defined function to compare Vader and Roberta Model through dist plot
341 def compModelsGraph(df):
342     plt.figure(figsize=(15, 8))
343     plt.suptitle("In-Depth VADER and ROBERTA Model Analysis", color="blue", fontsize=30)
344
345     plt.subplot(1, 3, 1)
346     plot_one = sns.distplot(list(df['vader_pos']), kde_kws=dict(linewidth=7))
347     plot_two = sns.distplot(df['roberta_pos'], kde_kws=dict(linewidth=7))
348     plt.xlabel("POSITIVE SCORES", color='red', fontsize=16)
349     plt.ylabel("DENSITY", color='red', fontsize=16)
350     plt.legend(["VADER MODEL", "ROBERTA MODEL"])
351
352     plt.subplot(1, 3, 2)
353     plot_one = sns.distplot(list(df['vader_neg']), kde_kws=dict(linewidth=7))
354     plot_two = sns.distplot(df['roberta_neg'], kde_kws=dict(linewidth=7))
355     plt.xlabel("NEGATIVE SCORES", color='red', fontsize=16)
356     plt.ylabel("DENSITY", color='red', fontsize=16)
357     plt.legend(["VADER MODEL", "ROBERTA MODEL"])
358
359     plt.subplot(1, 3, 3)
360     plot_one = sns.distplot(list(df['vader_comp']), kde_kws=dict(linewidth=7))
361     plot_two = sns.distplot(df['roberta_neu'], kde_kws=dict(linewidth=7))
362     plt.xlabel("COMPOUND SCORES", color='red', fontsize=16)
363     plt.ylabel("DENSITY", color='red', fontsize=16)
364     plt.legend(["VADER MODEL", "ROBERTA MODEL"])
365
366     return plt.show()
367
368
369 # Plotting Vader Scores vs Ratings
370 def plotVaderResults(df):
371     fig, axs = plt.subplots(1, 3, figsize=(16, 8))
372     sns.barplot(data=df, x='review_ratings', y='vader_comp', ax=axs[0], ci=None)
373     axs[0].set_title("Compound Score (VADER)")
374
375     sns.barplot(data=df, x='review_ratings', y='vader_neg', ax=axs[1], ci=None)
376     axs[1].set_title("Negative Score (VADER)")
377
378     sns.barplot(data=df, x='review_ratings', y='vader_pos', ax=axs[2], ci=None)
379     axs[2].set_title("Positive Score (VADER)")
380
381     plt.tight_layout()
382
383     return plt.show()
384
385
386 # Plotting Roberta Model Scores vs Ratings
387 def plotRobertaResults(df):
388     fig, axs = plt.subplots(1, 3, figsize=(16, 8))
389     sns.barplot(data=df, x=df['review_ratings'], y=df['roberta_neu'], ax=axs[0], ci=None)
390     axs[0].set_title("Neutral Score (ROBERTA MODEL)")
391
392     sns.barplot(data=df, x=df['review_ratings'], y=df['roberta_neg'], ax=axs[1], ci=None)
393     axs[1].set_title("Negative Score (ROBERTA MODEL)")

```

```

393
394     sns.barplot(data=df, x=df['review_ratings'], y=df['roberta_pos'], ax= axs[2], ci=None)
395     axs[2].set_title("Positive Score (ROBERTA MODEL)")
396
397     plt.tight_layout()
398     return plt.show()
399
400
401 # Density Plot for final sentiment density (Date wise)
402 def sentimentDensityPlot(df):
403     plt.figure(figsize=(12, 7))
404     sns.kdeplot(data=df, x='review_date', hue="sentiment", palette="icefire",
405                  multiple="fill")
406     plt.title("Weighted Sentiment Spread", fontsize=30, color='b')
407     plt.xlabel("Review Date", fontsize=16, color="red")
408     plt.ylabel("Density", fontsize=16, color="red")
409     return plt.show()
410
411
412 # Bar Chart to compare the Sentiment values of ROBERTA and VADER model
413 def modelValsComp(df):
414     rob_vad_vals = list(df["vader_sentiment"]) + list(df["roberta_sentiment"])
415     which_vals = ["VADER" for x in range(len(df))] + ["ROBERTA" for x in range(len(df))]
416     temp_df = pd.DataFrame({"Sentiment": rob_vad_vals, "Model": which_vals})
417
418     plt.figure(figsize=(12, 7))
419     sns.histplot(data=temp_df, x="Sentiment", hue="Model", multiple="dodge", shrink=.8)
420     plt.title("ROBERTA VS VADER MODEL SENTIMENT COMPARISON", fontsize=30, color='b')
421     plt.xlabel("Sentiment", fontsize=16, color="red")
422     plt.ylabel("Review Count", fontsize=16, color="red")
423     return plt.show()
424
425
426 # Bar Chart and Pie Chart of weighted Sentiment Values
427 def sentimentPlot(df):
428     pos_count = sum(df["sentiment"] == "Positive")
429     neu_count = sum(df["sentiment"] == "Neutral")
430     neg_count = sum(df["sentiment"] == "Negative")
431
432     print("-> Positive Sentiment Reviews : ", pos_count)
433     print("-> Neutral Sentiment Reviews : ", neu_count)
434     print("-> Negative Sentiment Reviews : ", neg_count)
435
436     plt.figure(figsize=(12, 7))
437     plt.rcParams["axes.edgecolor"] = "black"
438     plt.rcParams["axes.linewidth"] = 2.50
439     plt.suptitle("Reviews Sentiment Chart (Weighted 3:2 (ROBERTA MODEL : VADER MODEL))", fontsize=30, color
440                 ='b')
441     plt.subplot(1, 2, 1)
442
443     plt.bar(x=["Positive", "Neutral", "Negative"],
444             height=[pos_count, neu_count, neg_count],
445             color=['green', 'yellow', 'red'])
446
447     plt.xlabel("Sentiment", fontsize=16, color="red")
448     plt.ylabel("Count of reviews", fontsize=16, color="red")
449
450     plt.subplot(1, 2, 2)
451     plt.pie(x=[pos_count, neu_count, neg_count], labels=["Positive", "Neutral", "Negative"],
452             colors=['green', 'yellow', 'red'], autopct='%.0f%')
453
454     plt.tight_layout()
455     return plt.show()
456
457 """
458 MAIN FUNCTION
459 """
460
461
462 def mainReviewsAnalysis(df, product_name, cmts_size):
463     # Creating a subset of only Amazon related products
464     amz_df = df[df["ecommerce_website"] == "Amazon"]
465
466     # Scrapping the amazon product's review using the module created by me in Webscrapping.py
467     comments_scraping = ScrapeComments(amz_df, product_name, cmts_size)
468     df = comments_scraping.getDf()
469
470     # Checking once more for null/na/redundant values in the freshly scraped df

```

```
471     if 'Unnamed: 0' in df.columns:
472         df.drop(['Unnamed: 0'], axis=1, inplace=True)
473
474     cleanCommentsDf(df)
475     addReviewLen(df)
476     # Using the user defined function to generate vader and roberta model values for the comments
477     vader_df = getVaderDf(df)
478     roberta_df = getRobertaDf(df)
479
480     # Joining the initial data
481     x = roberta_df.join(vader_df)
482     df_cmt = df.join(x)
483     getVaderSentiment(df_cmt)
484     getRobertaSentiment(df_cmt)
485     getFinalSentiment(df_cmt)
486
487     return df_cmt
488
489
490 # Visualising Plots function with the help of the otehr methods created
491 def visualiseCommentsAnalysis(df_cmts):
492     plt0 = show_wordcloud(list(df_cmts['review_content']))
493     plt1 = sentimentPlot(df_cmts)
494     plt2 = sentimentDensityPlot(df_cmts)
495     plt3 = modelValsComp(df_cmts)
496     plt4 = ratingsItemPlot(df_cmts)
497     plt5 = sizeCompPlot(df_cmts)
498     plt6 = compModelsGraph(df_cmts)
499     plt7 = plotVaderResults(df_cmts)
500     plt8 = plotRobertaResults(df_cmts)
501
```