

“MACHINE LEARNING FOR E-COMMERCE PRODUCT EVALUATION,RANKING AND RECOMMENDATION WITH SENTIMENTAL ANALYSIS”

by

S. M. AKSHAYA 2021115012

S. SATHYA DHARINI 2021115097

*A summer project report
submitted to the Faculty of*

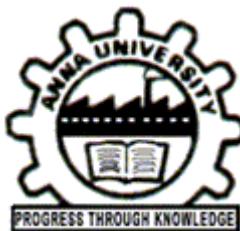
INFORMATION AND COMMUNICATION ENGINEERING

*for the partial fulfillment of the
award of the degree of*

Bachelor of Technology

in

Information Technology



**DEPARTMENT OF INFORMATION SCIENCE AND TECHNOLOGY
COLLEGE OF ENGINEERING, ANNA UNIVERSITY
CHENNAI 600 025.**

AUGUST 2023

CERTIFICATE

Certified that this report titled "**MACHINE LEARNING FOR E-COMMERCE PRODUCT EVALUATION ,RANKING AND RECOMMENDATION WITH SENTIMENTAL ANALYSIS**", is a *bona fide* work of **Ms. AKSHAYA S M (2021115012), Ms. SATHYA DHARINI S (2021115097)** who carried out the work under my supervision, for the partial fulfillment of the requirements for the award of the degree of *Bachelor of Engineering in Information Technology*. Certified further that to the best of my knowledge and belief, the work reported herein does not form part of any other thesis or dissertation on the basis of which a degree or an award was conferred on an earlier occasion.

Place: Chennai.

Date:

Dr. K INDRA GANDHI

Associate Professor,
Dept. of Information Science and Techn.,
College of Engineering, Guindy
Anna University,
Chennai - 600 025.

COUNTERSIGNED

Dr. S. Sridhar
Professor and Head
Department of Information Science and Technology
College of Engineering, Anna University,
Chennai 600 025.

ACKNOWLEDGEMENT

I would like to express my sincere thanks and deep sense of gratitude to my guide *Dr. K. INDRA GANDHI*, Associate Professor, Department of Information Science and Technology, for her valuable guidance suggestions and constant encouragement paved way for the successful completion of the project work.

I also thank *Dr. S. Sridhar*, Professor & Head, Department of Information Science and Technology for his kind support and for providing necessary facilities to carry out the work.

I am bound to convey my sincere thanks to the project committee *Dr. P. Varalakhsmi*, Professor, *Dr. S. Bama*, Associate Professor, and *Dr. K. Arul Deepa*, Assistant Professor, Department of Information Science and Technology for their suggestions and motivations.

S. M. AKSHAYA

S. SATHYADHARINI

ABSTRACT

The "Machine Learning for E-commerce Product Evaluation ,Ranking and Recommendation with Sentimental Analysis" project aims to simplify the process of finding the perfect product on an ecommerce site for users by leveraging a combination of machine learning (ML) algorithms and web scraping techniques. By scraping data from Amazon, including product descriptions, URLs, review counts, overall ratings, and specific rating percentages, the system gathers the necessary information. This data is then processed and analyzed using Python, using clustering algorithms for optimal pricing and sentiment analysis performed using Hugging Face's pretrained NLP model.

The primary goal of this automation is to save user's time and provide accurate product recommendations based on ratings, reviews, and customer needs(specifications). Instead of analysing through thousands of products datas on a ecommerce site, users can rely on this system to rank products based on various filters, including data deduplication, data preprocessing, price optimization through clustering, review emotion classification using Hugging Face's pretrained ML model, and sentiment analysis using Vader scores. The final output includes Vader scores along with the top-ranked product's URL, ensuring a streamlined and efficient product selection process.

TABLE OF CONTENTS

CHAPTER NO.	TITLE	PAGENO.
	ABSTRACT	4
	LIST OF ABBREVIATIONS	7
1	INTRODUCTION	8
	1.1 OVERVIEW	
	1.2 LITERATURE OF SURVEY	
	1.3 ORGANIZATION OF THE REPORT	11
2	REQUIREMENTS SPECIFICATION	12
	2.1 SOFTWARE SPECIFICATIONS	12
3	SYSTEM ORGANIZATION	13
	3.1 SYSTEM ARCHITECTURE DESIGN	13
4	IMPLEMENTATION AND RESULTS	17
	4.1 IMPLEMENTATION	17
	4.2 RESULTS	19
5	CONCLUSION AND FUTURE WORK	21
	5.1 CONCLUSION	21
	5.2 FUTUREWORK	21
	5.3 REFERENCES	22

CHAPTER 1

INTRODUCTION

1.1 OVERVIEW

Literature in the following areas has been analyzed:

- a. WEB SCRAPPING
- b. DATA PROCESSING USING PYTHON, CLASSIFICATION USING HUGGING PRE-TRAINED MODEL AND SENTIMENTAL ANALYSIS

WEB SCRAPING:-

Introduction

Web scraping is a technique used to extract data from websites for various purposes, including research, analysis, and business intelligence. This literature survey explores the historical context, legal and ethical aspects, techniques, applications, challenges, tools, and future trends in web scraping.

Historical Perspective

Web scraping has evolved alongside the internet. In the early days, it was often referred to as "screen scraping" and was used primarily for automating data entry tasks. Over time, web scraping has become more sophisticated and is now integral to data-driven decision-making processes.

Legal and Ethical Aspects

The legal landscape of web scraping is complex and varies by jurisdiction. Copyright, data protection, and terms of service agreements play a significant role. Ethical

considerations include issues of privacy, consent, and responsible data usage. Legal cases, such as HiQ v. LinkedIn, have shaped the legal understanding of web scraping.

Web Scraping Techniques

Web scraping techniques range from simple HTML parsing to more advanced methods like web automation and API integration. Tools such as BeautifulSoup, Scrapy, and Selenium are commonly used to facilitate scraping. The choice of technique depends on the specific requirements of the task.

Applications of Web Scraping

Web scraping finds applications across industries, including e-commerce, finance, healthcare, and academia. Examples include price monitoring in e-commerce, sentiment analysis in social media, and data collection for academic research.

References used:

<https://www.selenium.dev/about/>

<https://pypi.org/project/beautifulsoup4/>

<https://chromedriver.chromium.org/downloads>

DATA PROCESSING AND SENTIMENTAL ANALYSIS:-

Introduction

Pretrained models have revolutionized the field of Natural Language Processing (NLP) by providing powerful language representations that can be fine-tuned for various downstream tasks. This literature survey explores the historical context, architecture, training techniques, applications, challenges, and future trends in pretrained models for NLP.

Applications of Pretrained Models

Pretrained models have been applied to a wide range of NLP tasks, including sentiment analysis, named entity recognition, machine translation, question-answering, and text generation. They serve as strong feature extractors and can be fine-tuned with relatively small task-specific datasets.

Tools and Frameworks

Frameworks like Hugging Face Transformers have made it easier for researchers and developers to work with pretrained models. They offer prebuilt models and APIs for fine-tuning and inference, facilitating their adoption in various applications.

Future Trends

The future of pretrained models in NLP is promising. Ongoing research focuses on reducing model size and inference latency, improving interpretability, and addressing ethical concerns. Multilingual models, few-shot learning, and zero-shot learning are emerging trends.

References used:

<https://huggingface.co/>

<https://www.kaggle.com/>

1.2 LITERATURE SURVEY

1. "EFFICIENT SCRAPING OF DATA FROM WEBSITES USING SELENIUM"

This paper concludes about the benefits of using selenium as web scrapper which reduces the manual work of extracting data from different websites and thus the automatic generation of data makes easier to process and analyse further.

Authors: [1]Shreya V. Dhone student ,[2] Anupama D. Sakhare Asst Prof,[3] Satish J.Sharma Prof.

1,2,3 Dept of Electronics and Computer Science, Rakshtrasant Tukdoji Maharaj Nagpur University,India . 2022 jetir, volume 9, issue 6,issn:2349-5162

2."RoBERTa-GRU: A Hybrid Deep Learning Model for Enhanced Sentiment Analysis"

The results of the experiments carried out on this paper IMDb, Sentiment140, and Twitter US Airline datasets indicate that the RoBERTa-GRU model outperforms all other comparison methods, with accuracy scores of 94.63%, 89.59%, and 91.52%, respec-

tively. The combination of the RoBERTa and GRU results in a powerful and efficient model for sentiment analysis, making it a promising solution for a variety of NLP tasks.

Authors: Kian Long Tan , Chin Poo Lee, Kian Ming Lim , Multimedia University, Malaysia

3. "The k-means Algorithm: A Comprehensive Survey and Performance Evaluation"

The paper made experimental analysis evaluated the performance of k means clustering algorithm, and proposing where it serves the best results i.e differing and being dependent on the nature dataset.

Authors: Mohiuddin Ahmed, Raihan Seraj, Syed Mohammed Shamsul Islam

4. "Sentiment analysis of product reviews: A review "

The paper made experimental analysis by doing sentimental analysis on various products reviews and proposes Support Vector Machine giving the high accuracy.

Author: T.K. Shivaprasad, Jyothi Shetty, NMAM Institute of Technology

Published in: 2017 International Conference on Inventive Communication and Computational Technologies (ICICCT)

5. "A review on sentiment analysis and emotion detection from text"

This paper concludes that the lexicon based technique performs well in both sentiment and emotional analysis.

Authors: Nandwani, Rupali Verma

6. "VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text"

This paper concludes how much is good to use the vader scores for sentimental analysis

Authors: C.J. Hutto, Eric Gilbert , Georgia Institute Of Technology

7. "An overview and comparison of free Python libraries for data mining and big data analysis"

This paper says about how much python is best language in the data -analysis and machine learning

Authors: Stacin, Jovic

1.3 ORGANIZATION OF THE REPORT

The rest of the report is organized as follows.

Chapter 2 describes the software requirements specifications.

Chapter 3 describes and shows the overall Architectural design of the system that has been developed.

Chapter 4 provides the results obtained by the system and the performance analysis of the system.

Chapter 5 deals with the concluding remarks about the project and the future areas of development in the project.

CHAPTER 2

SOFTWARE REQUIREMENTS SPECIFICATION

2.1 SOFTWARE SPECIFICATION

- ❖ **Python:** Code is written in python. Compatible version (3.x)
- ❖ **Bs4:** Beautiful Soup. Used for HTML parsing.
- ❖ **Selenium:** Used for browser automation.
- ❖ **Web Driver (Chrome):** Used to automate browser automations for the Google Chrome web browser

❖ LIBRARIES INCLUDED:

- 1)csv
- 2)selenium
- 3)pandas
- 4)re
- 5)matplotlib.pyplot
- 6)sklearn.cluster
- 7)os
- 8)transformers
- 9)nltk
- 10)libraries
- 11)strings

CHAPTER 3

SYSTEM DESIGN

3.1 SYSTEM ARCHITECTURE:-

a. Data Scraping:

It includes a web scraper for Amazon product search results. It takes a user-entered product as input, uses BeautifulSoup and Selenium to scrape data from Amazon's search results pages, and then saves the extracted information, including product descriptions, prices, ratings, review counts, and URLs, into a CSV file named 'font67.csv'. The script iterates through multiple pages of search results and ensures that the extracted data is structured and saved for analysis or other purposes. It's important to use such web scraping scripts responsibly and in compliance with the website's terms of service and legal requirements.

b. User specifications:

It loads data from a CSV file named 'font67.csv' and filters out rows with missing values in specified columns like 'Description', 'Price', 'rating', 'reviewcount', and 'url'. It allows the user to input multiple keywords for product descriptions and filters the data based on these keywords.

c. Optimum price range calculation:(K MEANS CLUSTERING)

It calculates the distribution of review counts and divides the data into ranges. It performs clustering on the 'Price' column to identify clusters with the most data points. i.e., which are classified by k-means clustering algorithm considering products with high review counts and grouping them into ranges and considering the price ranges where most of the products lie in ,and filtering the datas in a separate csv file and proceeding with the next process for filtering out.

d. Scraping rating's percentage and filtering:

It uses Selenium to scrape additional data from Amazon product pages, including rating percentages for 5-star and 4-star ratings. It removes duplicate rows from the data based on the 'Description' column and saves the deduplicated data to 'filtered_data9_deduplicated.csv'. It combines the scraped data with the original data, including rating percentages and other information. It computes a combined score for each row based on selected columns like '5star', '4star', 'reviewcount', and 'o_rating'. It sorts the data based on the combined score in descending order. It saves the sorted data to a new CSV file named 'sorted_data1.csv'.

e. Reviews Scraping:

The script defines functions to scrape regular reviews and critical reviews from Amazon product pages.

It extracts review text by navigating to review pages, scraping the reviews, and handling pagination. The script loads product data from a CSV file ('sorted_data1.csv') containing product descriptions and URLs.

It creates a list ('products_data') to store product information. It iterates through the product data and for each product, collects its reviews (both regular and critical).

The reviews are written to separate CSV files ('merix1.csv' for regular reviews and 'scrapzz.csv' for critical reviews).

f. Sentimental Analysis:

The Hugging Face Transformers library is used to perform sentiment analysis on the scraped reviews.

A pre-trained model ("SamLowe roberta-base-go_emotions") is employed to classify the sentiments within the text. It analyses and counts the sentiment labels present in the text data from the CSV file, providing insights into the distribution of sentiments within the text.

g. Vader Score Calculation:

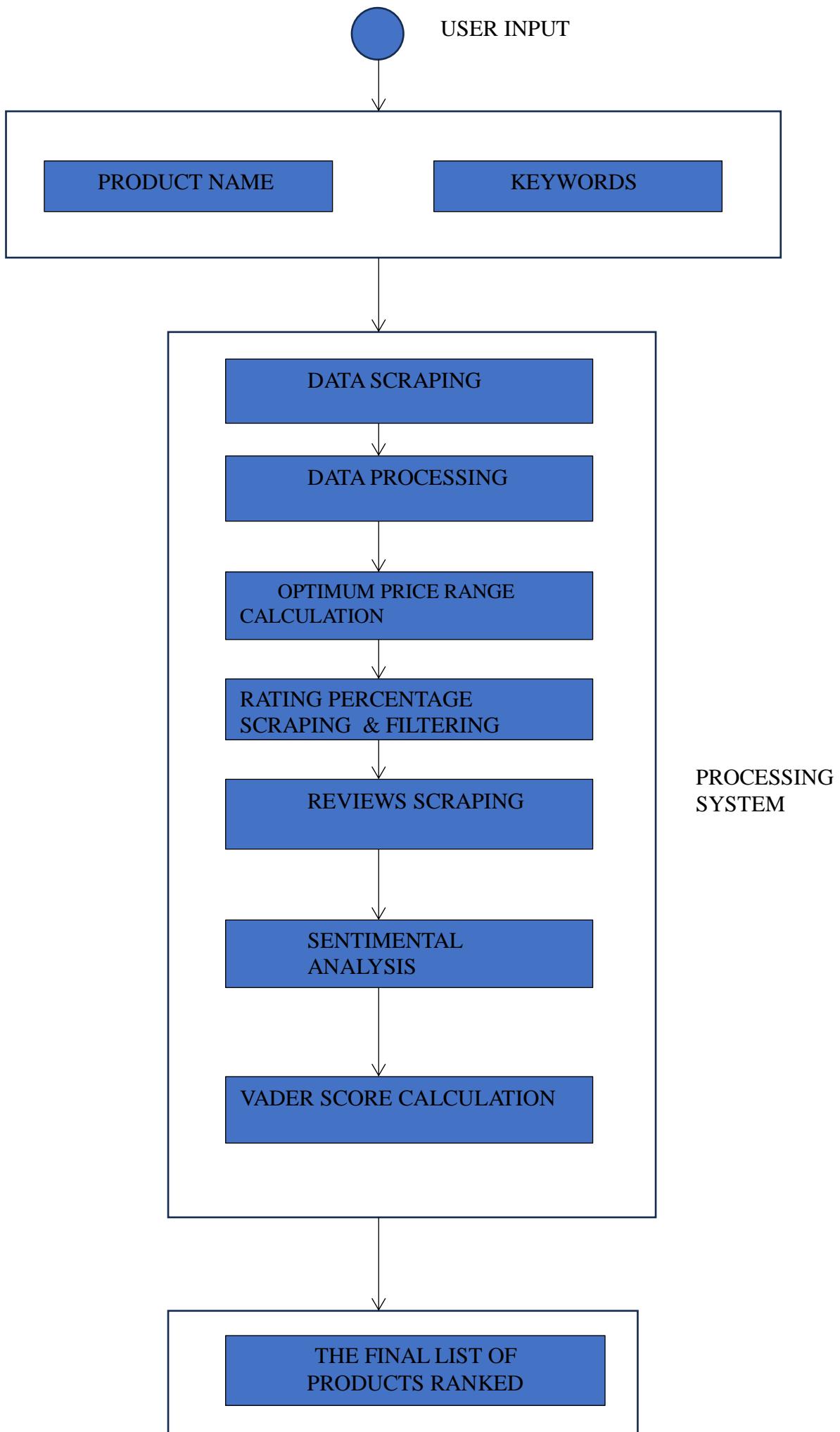
NLTK resources such as a tokenizer, stopwords, a lemmatizer, and the VADER sentiment intensity analyzer are initialized.

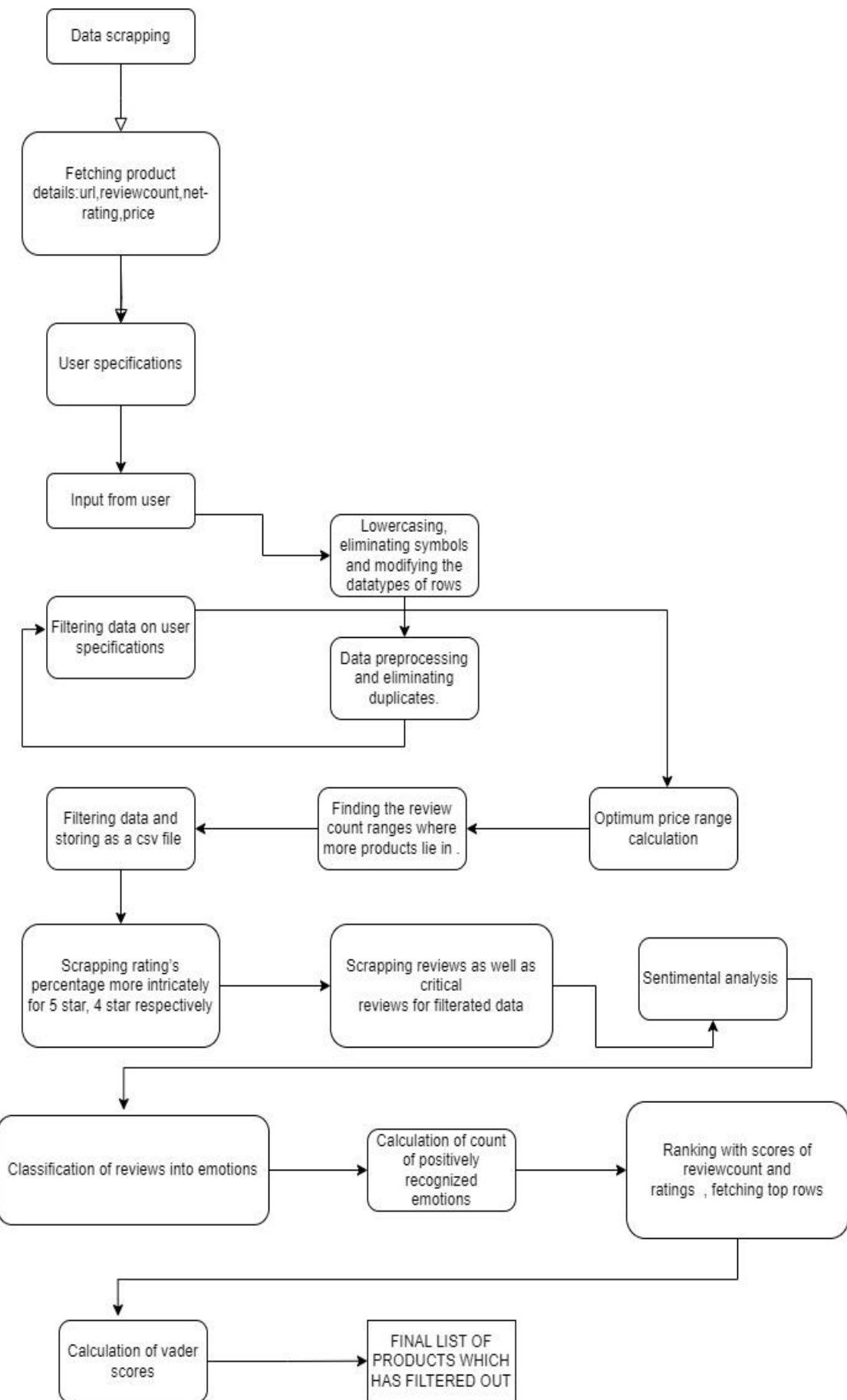
The code reads a CSV file ('merix1.csv') and preprocesses the text in each cell. Sentiment analysis using VADER is performed on the preprocessed words, and the compound scores are summed for each cell.

The summed sentiment scores are printed for each row.

h. Final Output:

Based on VADER score, final output is displayed. The product with higher VADER Score is the best recommended and considered by a lot of people as Amazon users.



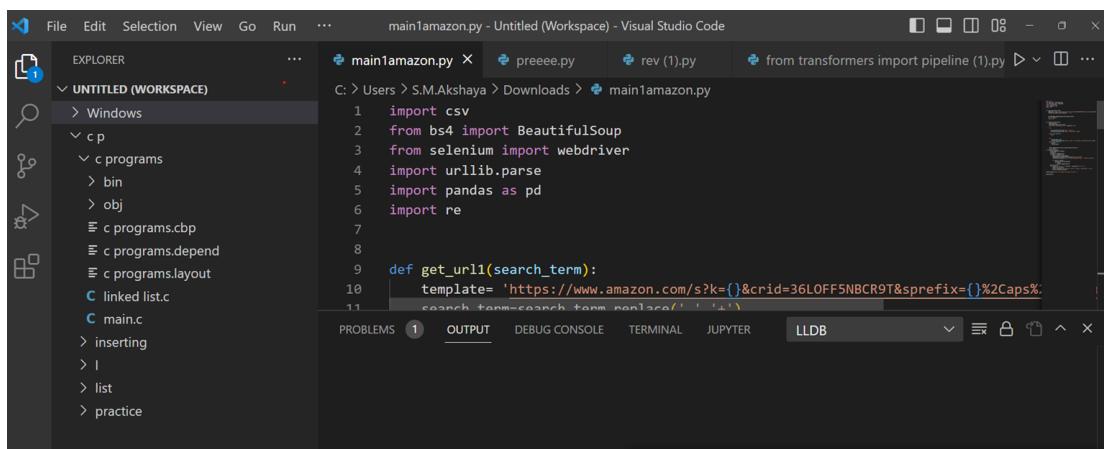


IMPLEMENTATION AND RESULTS

4.1 IMPLEMENTATION

4.1.1 Web Scraping:

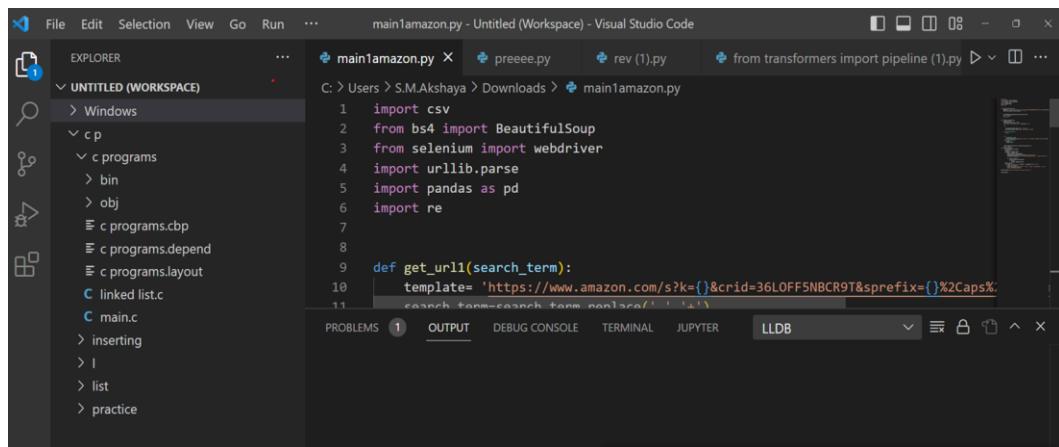
a.Importing Libraries:



```
File Edit Selection View Go Run ... main1amazon.py - Untitled (Workspace) - Visual Studio Code
EXPLORER File Explorer
UNTITLED (WORKSPACE)
> Windows
> c p
> c programs
> bin
> obj
> c programs.cbp
> c programs.depend
> c programs.layout
C linked list.c
C main.c
> inserting
> i
> list
> practice
main1amazon.py X preeee.py rev (1).py from transformers import pipeline (1).py ...
C: > Users > S.M.Akshaya > Downloads > main1amazon.py
1 import csv
2 from bs4 import BeautifulSoup
3 from selenium import webdriver
4 import urllib.parse
5 import pandas as pd
6 import re
7
8
9 def get_url1(search_term):
10     template= 'https://www.amazon.com/s?k={}&crid=36LOFF5NBCR9T&sprefix={}%2Caps%2B'
11     search_term=search_term.replace(' ','+')
PROBLEMS 1 OUTPUT DEBUG CONSOLE TERMINAL JUPYTER LLDB
```

Importing all the necessary libraries including: csv for saving the data in excel, BeautifulSoup for parsing a web page and Selenium for web automation and a driver has been installed for the Google chrome browser and pandas for data manipulation and analysis and finally importing re ‘regular expressions’ modules.

b. Define Function “`get_url1(search_term)`”

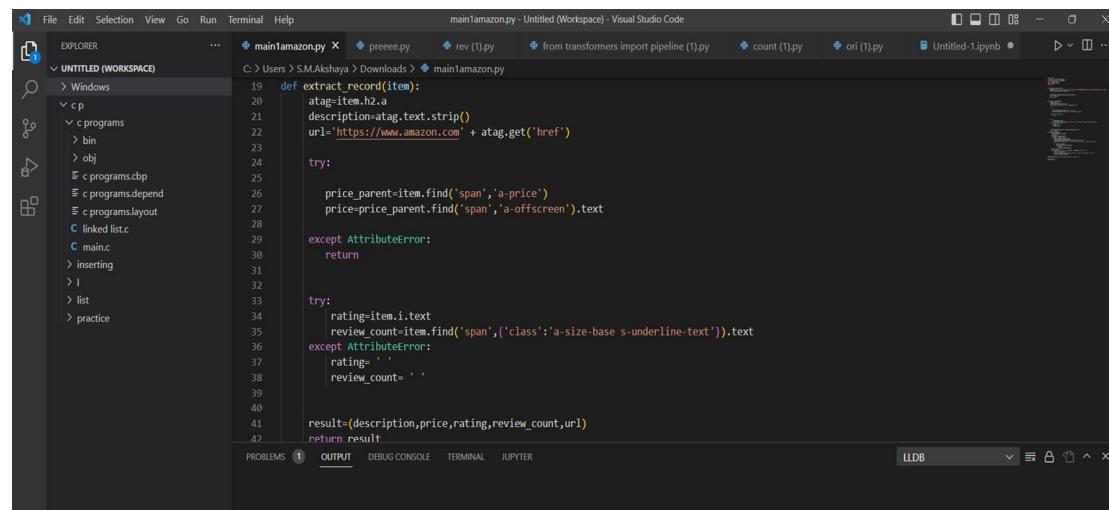


```
File Edit Selection View Go Run ... main1amazon.py - Untitled (Workspace) - Visual Studio Code
EXPLORER main1amazon.py x preeee.py rev (1).py from transformers import pipeline (1).py ...
UNTITLED (WORKSPACE)
C: > Users > S.M.Akshaya > Downloads > main1amazon.py
1 import csv
2 from bs4 import BeautifulSoup
3 from selenium import webdriver
4 import urllib.parse
5 import pandas as pd
6 import re
7
8
9 def get_url1(search_term):
10     template= 'https://www.amazon.com/s?k={} & crid=36LOFF5NBCR9T&sprefix={} %2Caps%2C'
11     search_term=search_term.replace(' ','+')
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
```

The screenshot shows the Visual Studio Code interface with the file `main1amazon.py` open. The code defines a function `get_url1` that takes a `search_term` and returns a URL template. The template includes the Amazon search URL with placeholders for the search term and category ID.

Defining a user-defined function named `get_url1(search_term)` so this helps to process the target url, Where the user specifies the product name, thus we add that as a substring to the common url and hence we get the target url which helps to direct the appropriate site containing the product details.

c. Define function “`extract_record(item)`”:

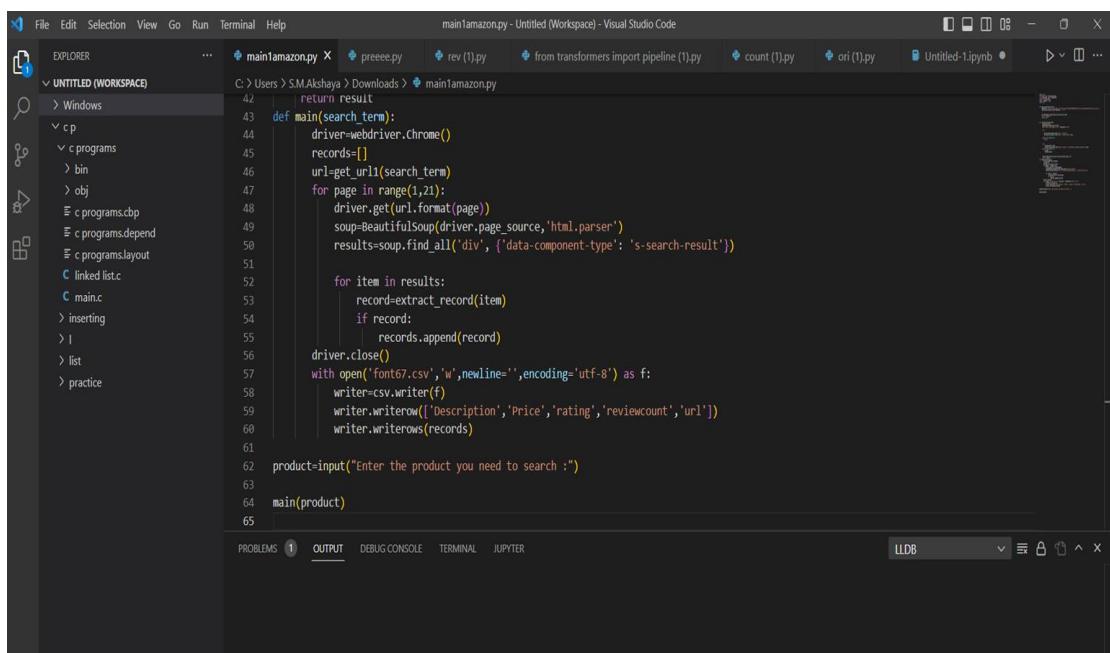


```
File Edit Selection View Go Run Terminal Help main1amazon.py - Untitled (Workspace) - Visual Studio Code
EXPLORER main1amazon.py x preeee.py rev (1).py from transformers import pipeline (1).py ... count (1).py ori (1).py Untitled-1.ipynb ...
UNTITLED (WORKSPACE)
C: > Users > S.M.Akshaya > Downloads > main1amazon.py
19 def extract_record(item):
20     atag=item.h2.a
21     description=atag.text.strip()
22     url='https://www.amazon.com' + atag.get('href')
23
24     try:
25         price_parent=item.find('span','a-price')
26         price=price_parent.find('span','a-offscreen').text
27
28     except AttributeError:
29         return
30
31     try:
32         rating=item.i.text
33         review_count=item.find('span',{'class':'a-size-base s-underline-text'}).text
34
35     except AttributeError:
36         rating=''
37         review_count=''
38
39     result=(description,price,rating,review_count,url)
40
41     return result
42
```

The screenshot shows the Visual Studio Code interface with the file `main1amazon.py` open. The code defines a function `extract_record` that takes an item object and extracts its description, price, rating, review count, and URL.

This function used to scrap the datas :“description”, “price”, “url”, “review count” and the overall “rating” of the product.

d.Define the main function “**main(search_item)**”:

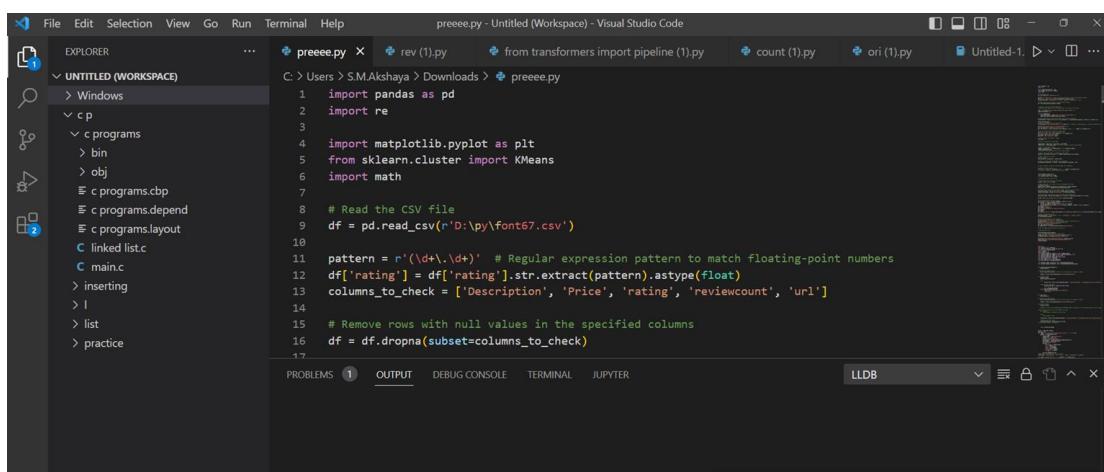


```
C:\Users\S.M.Akshaya>Downloads>mainAmazon.py
42     return result
43 def main(search_term):
44     driver=webdriver.Chrome()
45     records=[]
46     url=get_url1(search_term)
47     for page in range(1,21):
48         driver.get(url.format(page))
49         soup=BeautifulSoup(driver.page_source,'html.parser')
50         results=soup.find_all('div', {'data-component-type': 's-search-result'})
51
52         for item in results:
53             record=extract_record(item)
54             if record:
55                 records.append(record)
56
57         driver.close()
58         with open('font67.csv','w',newline='',encoding='utf-8') as f:
59             writer=csv.writer(f)
60             writer.writerow(['Description','Price','rating','reviewcount','url'])
61             writer.writerows(records)
62
63         product=input("Enter the product you need to search :")
64     main(product)
65
```

This function main almost scraps a maximum of 21 pages of details containing product details and storing all the datas in a excel sheet, this traverses to each and every page and scraps one by one.

4.1.2 Data Filtering and Price clustering: -

a.Data Loading and Preprocessing:



The screenshot shows the Visual Studio Code interface with the following details:

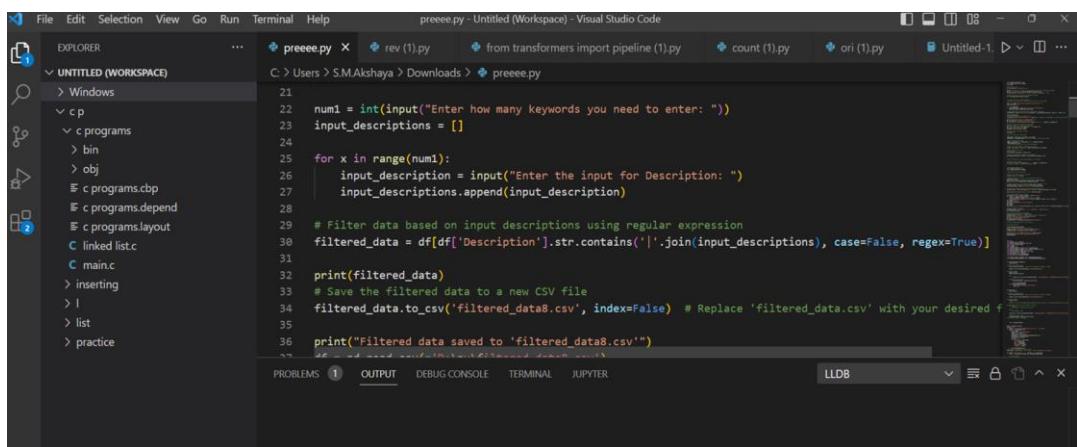
- File Explorer:** Shows a workspace named "UNTITLED (WORKSPACE)" containing files like "preeee.py", "rev (1).py", "from transformers import pipeline (1).py", "count (1).py", "ori (1).py", and "Untitled-1.py".
- Code Editor:** Displays Python code for reading a CSV file and performing data cleaning and conversion.

```
C: > Users > S.M.Akshaya > Downloads > preeee.py
1 import pandas as pd
2 import re
3
4 import matplotlib.pyplot as plt
5 from sklearn.cluster import KMeans
6 import math
7
8 # Read the CSV file
9 df = pd.read_csv("D:\\py\\font67.csv")
10
11 pattern = r'(\d+\.\d+)' # Regular expression pattern to match floating-point numbers
12 df['rating'] = df['rating'].str.extract(pattern).astype(float)
13 columns_to_check = ['Description', 'Price', 'rating', 'reviewcount', 'url']
14
15 # Remove rows with null values in the specified columns
16 df = df.dropna(subset=columns_to_check)
17
```

- Bottom Bar:** Includes tabs for PROBLEMS, OUTPUT, DEBUG CONSOLE, TERMINAL, and JUPYTER, along with LLDB and other icons.

Preprocessing the datas that has been already stored to “font67.csv” and converting them to the appropriate data types that helps in the further processing.

b.User Input for “FILTERING”:



The screenshot shows the Visual Studio Code interface with the following details:

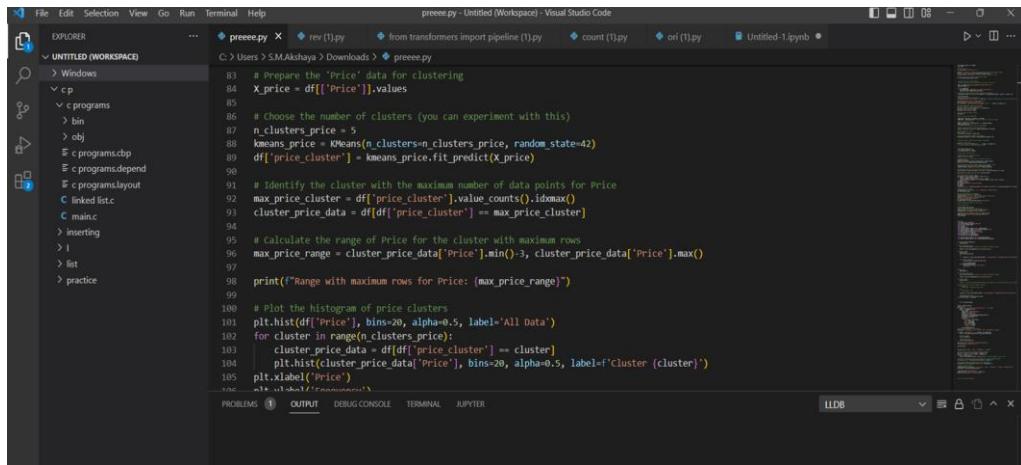
- File Explorer:** Shows a workspace named "UNTITLED (WORKSPACE)" containing files like "preeee.py", "rev (1).py", "from transformers import pipeline (1).py", "count (1).py", "ori (1).py", and "Untitled-1.py".
- Code Editor:** Displays Python code for taking user input for keywords and filtering a dataset based on those keywords.

```
21 num1 = int(input("Enter how many keywords you need to enter: "))
22 input_descriptions = []
23
24 for x in range(num1):
25     input_description = input("Enter the input for Description: ")
26     input_descriptions.append(input_description)
27
28 # Filter data based on input descriptions using regular expression
29 filtered_data = df[df['Description'].str.contains('|'.join(input_descriptions), case=False, regex=True)]
30
31 print(filtered_data)
32 # Save the filtered data to a new CSV file
33 filtered_data.to_csv('filtered_data8.csv', index=False) # Replace 'filtered_data.csv' with your desired file name
34
35 print("Filtered data saved to 'filtered_data8.csv'")
36
```

- Bottom Bar:** Includes tabs for PROBLEMS, OUTPUT, DEBUG CONSOLE, TERMINAL, and JUPYTER, along with LLDB and other icons.

This part asks the user to enter the number of keywords and the keywords also thus using these keywords , we check up in all the product's description which contain either of these keywords specified and filtering accordingly.

c.Price Clustering:

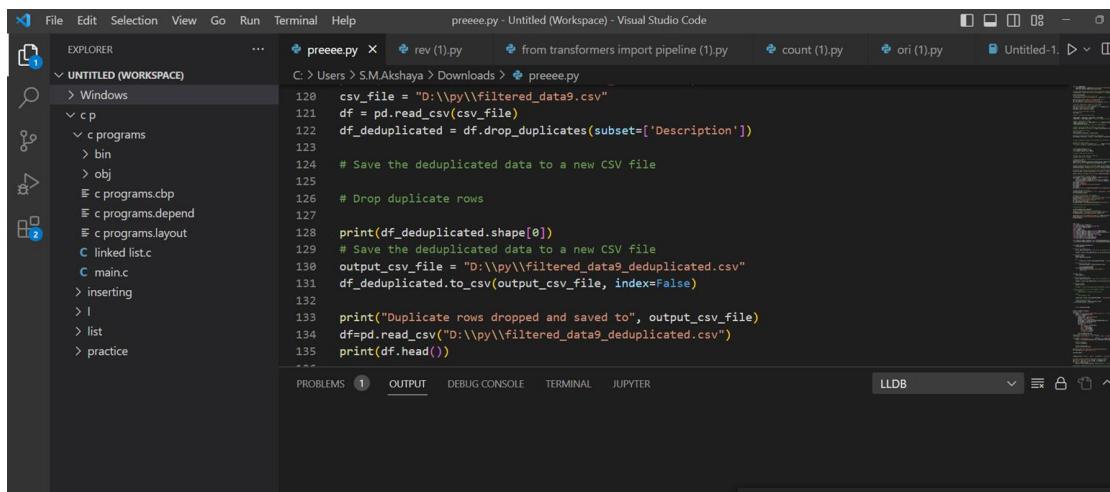


A screenshot of Visual Studio Code showing a Python script named 'preeee.py'. The code uses Scikit-Learn's KMeans algorithm to cluster data based on the 'Price' column. It then identifies the cluster with the maximum number of data points and calculates the price range for that cluster. The code includes comments explaining each step.

```
83 # Prepare the 'Price' data for clustering
84 x_price = df[['Price']].values
85
86 # choose the number of clusters (you can experiment with this)
87 n_clusters_price = 5
88 kmeans_price = KMeans(n_clusters=n_clusters_price, random_state=42)
89 df['price_cluster'] = kmeans_price.fit_predict(x_price)
90
91 # Identify the cluster with the maximum number of data points for Price
92 max_price_cluster = df['price_cluster'].value_counts().idxmax()
93 cluster_price_data = df[df['price_cluster'] == max_price_cluster]
94
95 # calculate the range of Price for the cluster with maximum rows
96 max_price_range = cluster_price_data['Price'].min(), cluster_price_data['Price'].max()
97
98 print(f"Range with maximum rows for Price: {max_price_range}")
99
100 # plot the histogram of price clusters
101 plt.hist(df["Price"], bins=20, alpha=0.5, label="All Data")
102 for cluster in range(n_clusters_price):
103     cluster_price_data = df[df['price_cluster'] == cluster]
104     plt.hist(cluster_price_data["Price"], bins=20, alpha=0.5, label=f"Cluster {cluster}")
105 plt.xlabel("Price")
106 plt.title("Histogram of Price Clusters")
```

First grouping them based on the review count and creating 5 ranges of rc's and counting the number as well as percentages of products falling to each range.Using the Scikit-Learn library's KMeans clustering algorithm to cluster data based on the 'Price' column into a specified number of clusters (n_clusters_price).Identifying the cluster with the maximum number of data points and calculating the price range for the cluster with the maximum rows.

d.Data deduplication:

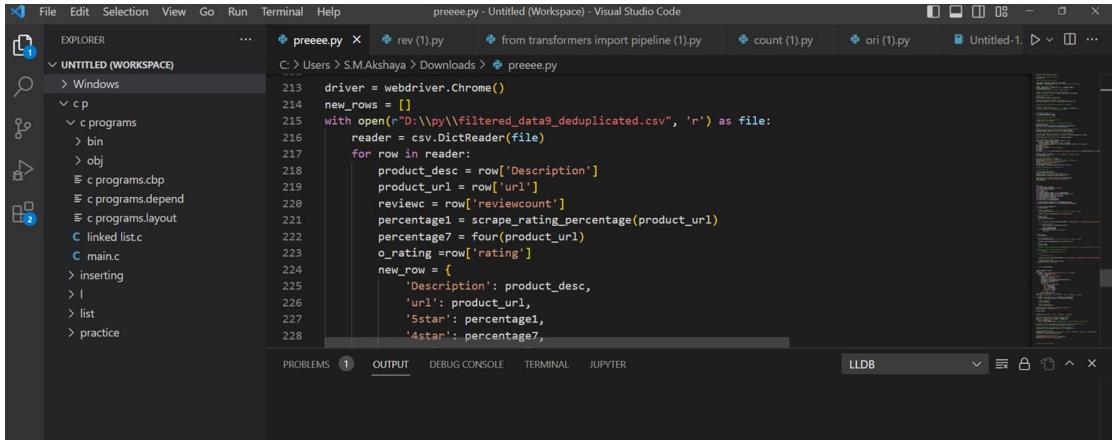


A screenshot of Visual Studio Code showing a Python script named 'preeee.py'. The code reads a CSV file, drops duplicate rows, and saves the deduplicated data to a new CSV file. The code includes comments explaining each step.

```
120 csv_file = "D:\\py\\\\filtered_data9.csv"
121 df = pd.read_csv(csv_file)
122 df_deduplicated = df.drop_duplicates(subset=['Description'])
123
124 # Save the deduplicated data to a new CSV file
125
126 # Drop duplicate rows
127
128 print(df_deduplicated.shape[0])
129 # Save the deduplicated data to a new CSV file
130 output_csv_file = "D:\\py\\\\filtered_data9_deduplicated.csv"
131 df_deduplicated.to_csv(output_csv_file, index=False)
132
133 print("Duplicate rows dropped and saved to", output_csv_file)
134 df=pd.read_csv("D:\\py\\\\filtered_data9_deduplicated.csv")
135 print(df.head())
```

Checking if any duplicate data rows are present and eliminating the duplicates

e.Scraping of rating percentages:

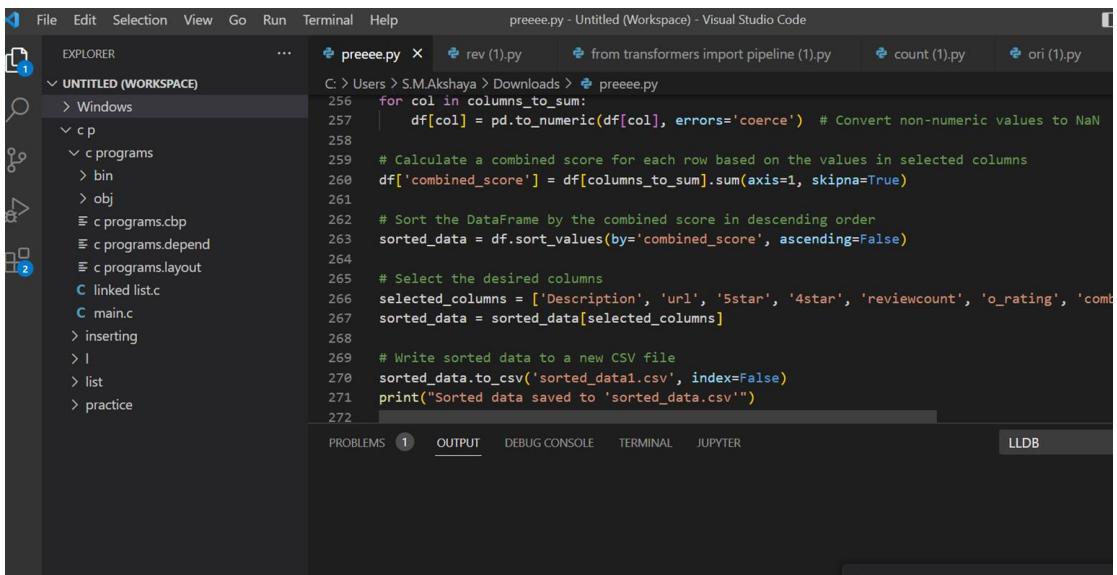


The screenshot shows the Visual Studio Code interface with the following details:

- File Bar:** File, Edit, Selection, View, Go, Run, Terminal, Help.
- Editor:** The main editor window displays a Python script named `preeee.py`. The code reads from a CSV file and calculates a percentage value based on the URL.
- Explorer:** On the left, the Explorer sidebar shows a workspace named "UNTITLED (WORKSPACE)" containing files like `c p`, `c programs`, and `main.c`.
- Bottom Navigation:** PROBLEMS, OUTPUT, DEBUG CONSOLE, TERMINAL, JUPYTER, LLDB.

Now scrapping the individual rating percentages including 5 star and 4 star and thus this helps in filtering the more product intricately.

f.Combining and Sorting data:



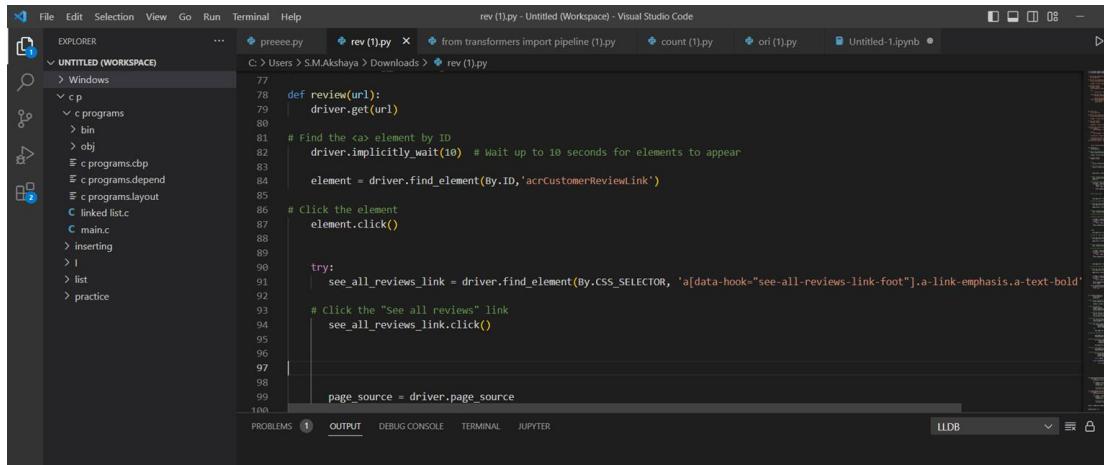
The screenshot shows the Visual Studio Code interface with the following details:

- File Bar:** File, Edit, Selection, View, Go, Run, Terminal, Help.
- Editor:** The main editor window displays a Python script named `preeee.py`. The code performs the following steps:
 - Converts non-numeric values to NaN.
 - Calculates a combined score for each row based on selected columns.
 - Sorts the DataFrame by the combined score in descending order.
 - Selects the desired columns.
 - Writes the sorted data to a new CSV file named `sorted_data1.csv`.
 - Prints a message indicating the data was saved to the new CSV file.
- Explorer:** On the left, the Explorer sidebar shows a workspace named "UNTITLED (WORKSPACE)" containing files like `c p`, `c programs`, and `main.c`.
- Bottom Navigation:** PROBLEMS, OUTPUT, DEBUG CONSOLE, TERMINAL, JUPYTER, LLDB.

Now finding the combined score of reviewcount , 5 star , 4 star and calculating the overall score(summation of three values) which helps in the ranking of the products hence , this process helps to filter the product and saving into new csv file .

4.1.3 Scraping Reviews:

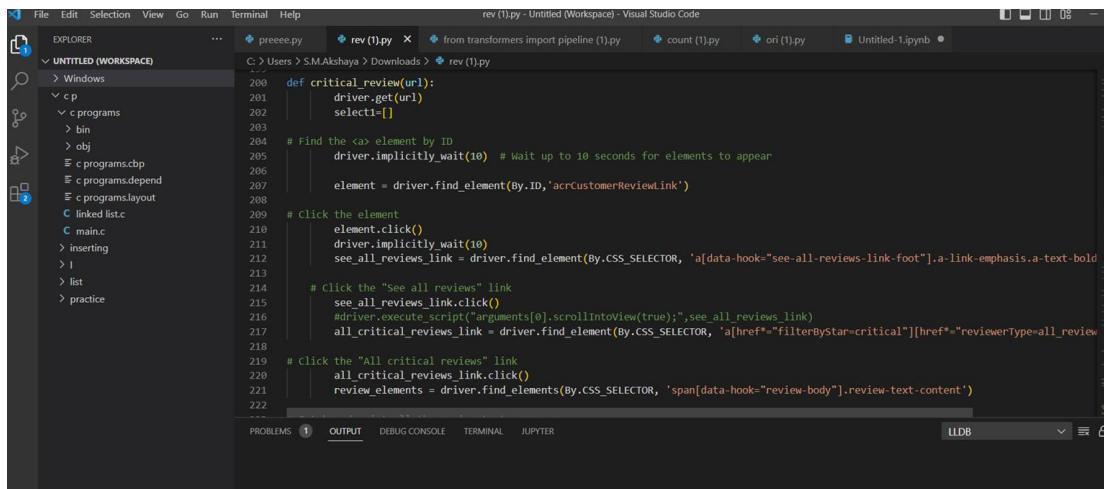
a. Web Scraping function:



The screenshot shows the Visual Studio Code interface with the following details:

- File Explorer:** Shows a workspace named "UNTITLED (WORKSPACE)" containing files like "rev (1).py", "preee.py", and "count (1).py".
- Code Editor:** Displays Python code for a "review" function. The code uses Selenium to navigate to a URL, find an element by ID, click it, and then click a "See all reviews" link to get the page source.
- Bottom Bar:** Shows tabs for PROBLEMS, OUTPUT, DEBUG CONSOLE, TERMINAL, and JUPYTER. The OUTPUT tab is selected.
- Right Side:** Includes a sidebar with file navigation and a status bar at the bottom right.

Now scrapping the reviews of the filtered data and this almost 50 -60 reviews of each product and get stored on a csv file.



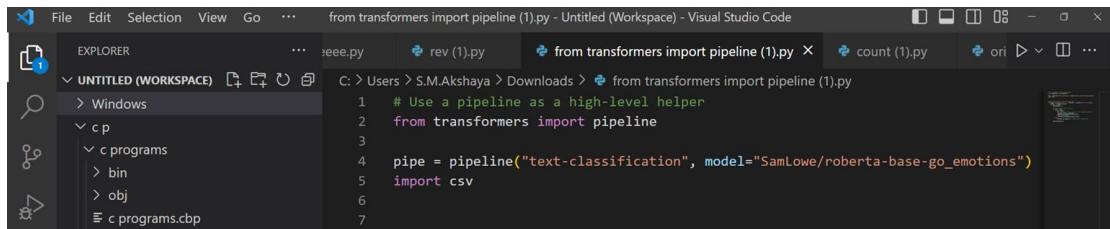
The screenshot shows the Visual Studio Code interface with the following details:

- File Explorer:** Shows a workspace named "UNTITLED (WORKSPACE)" containing files like "rev (1).py", "preee.py", and "count (1).py".
- Code Editor:** Displays Python code for a "critical_review" function. It follows a similar structure to the "review" function but includes additional logic to scroll into view and execute a script to filter reviews by star rating.
- Bottom Bar:** Shows tabs for PROBLEMS, OUTPUT, DEBUG CONSOLE, TERMINAL, and JUPYTER. The OUTPUT tab is selected.
- Right Side:** Includes a sidebar with file navigation and a status bar at the bottom right.

This helps in the scrapping of reviews which are categorized under a tag named "critical reviews"

4.1.4 Sentimental Analysis:

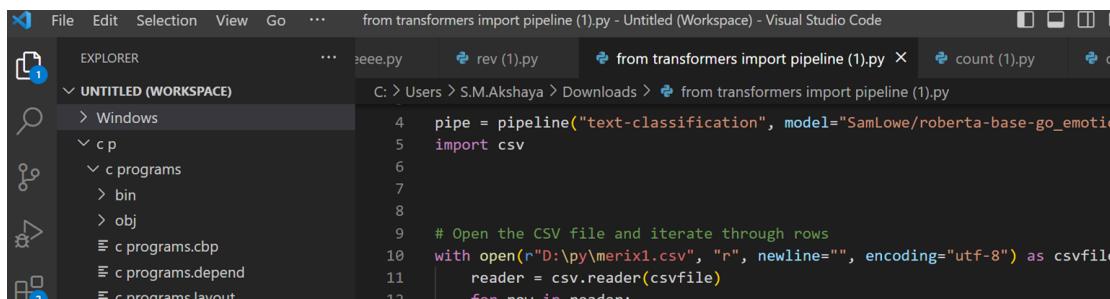
a.Pipeline Creation:



```
File Edit Selection View Go ... from transformers import pipeline (1).py - Untitled (Workspace) - Visual Studio Code
EXPLORER ... jeee.py rev (1).py from transformers import pipeline (1).py count (1).py ori (1).py ...
UNTITLED (WORKSPACE) > Windows > c p > bin > obj > c programs > c programs.cbp > c programs.depend > c programs.layout
C: > Users > S.M.Akshaya > Downloads > from transformers import pipeline (1).py
1 # Use a pipeline as a high-level helper
2 from transformers import pipeline
3
4 pipe = pipeline("text-classification", model="SamLowe/roberta-base-go_emotions")
5 import csv
6
7
```

Pretrained Roberta model has been used which has linked using application interface using hugging face , thus which has a ability to categorize the reviews to several emotion tag namely joy, sad, disappointment and more on,,,

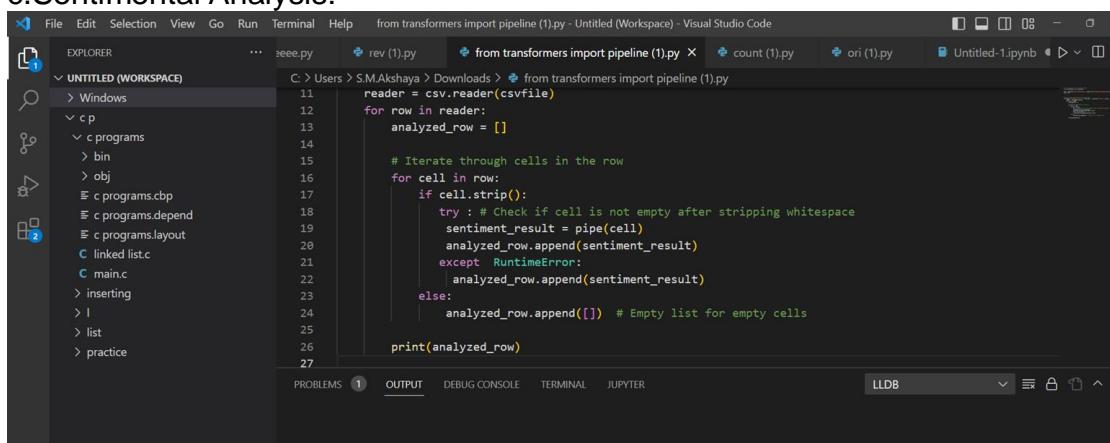
b.Csv Processing:



```
File Edit Selection View Go ... from transformers import pipeline (1).py - Untitled (Workspace) - Visual Studio Code
EXPLORER ... jeee.py rev (1).py from transformers import pipeline (1).py count (1).py ori (1).py ...
UNTITLED (WORKSPACE) > Windows > c p > bin > obj > c programs > c programs.cbp > c programs.depend > c programs.layout
C: > Users > S.M.Akshaya > Downloads > from transformers import pipeline (1).py
4 pipe = pipeline("text-classification", model="SamLowe/roberta-base-go_emotions")
5 import csv
6
7
8
9 # Open the CSV file and iterate through rows
10 with open(r"D:\py\merix1.csv", "r", newline="", encoding="utf-8") as csvfile
11     reader = csv.reader(csvfile)
12     for row in reader:
13         print(row)
```

Saving the reviews in a csv file .

c.Sentimental Analysis:

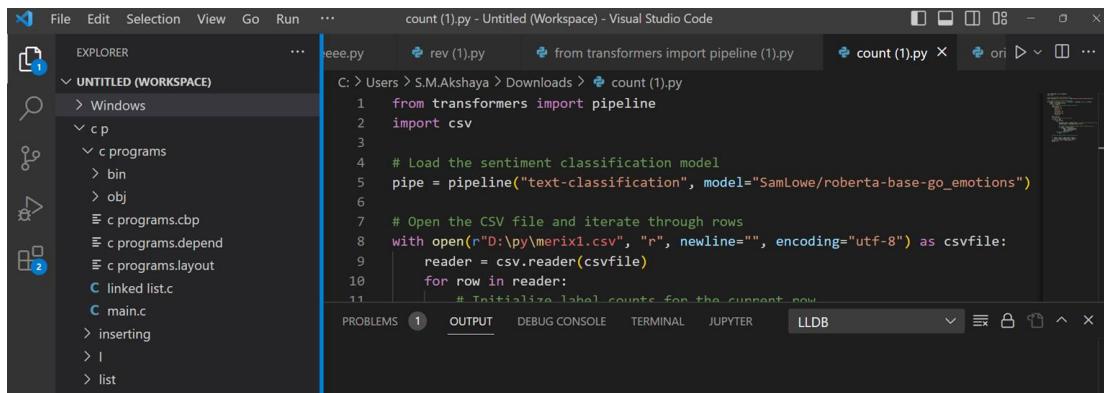


```
File Edit Selection View Go Run Terminal Help from transformers import pipeline (1).py - Untitled (Workspace) - Visual Studio Code
EXPLORER ... jeee.py rev (1).py from transformers import pipeline (1).py count (1).py ori (1).py Untitled-1.ipynb > ...
UNTITLED (WORKSPACE) > Windows > c p > bin > obj > c programs > c programs.cbp > c programs.depend > c programs.layout > inserting > I > list > practice
C: > Users > S.M.Akshaya > Downloads > from transformers import pipeline (1).py
11 reader = csv.reader(csvfile)
12 for row in reader:
13     analyzed_row = []
14
15     # Iterate through cells in the row
16     for cell in row:
17         if cell.strip():
18             try: # Check if cell is not empty after stripping whitespace
19                 sentiment_result = pipe(cell)
20                 analyzed_row.append(sentiment_result)
21             except RuntimeError:
22                 analyzed_row.append(sentiment_result)
23             else:
24                 analyzed_row.append([]) # Empty list for empty cells
25
26     print(analyzed_row)
27
```

Each and every review will be categorized into a emotion and labelled

4.1.5 Label Counting of reviews using a pre-trained Model:

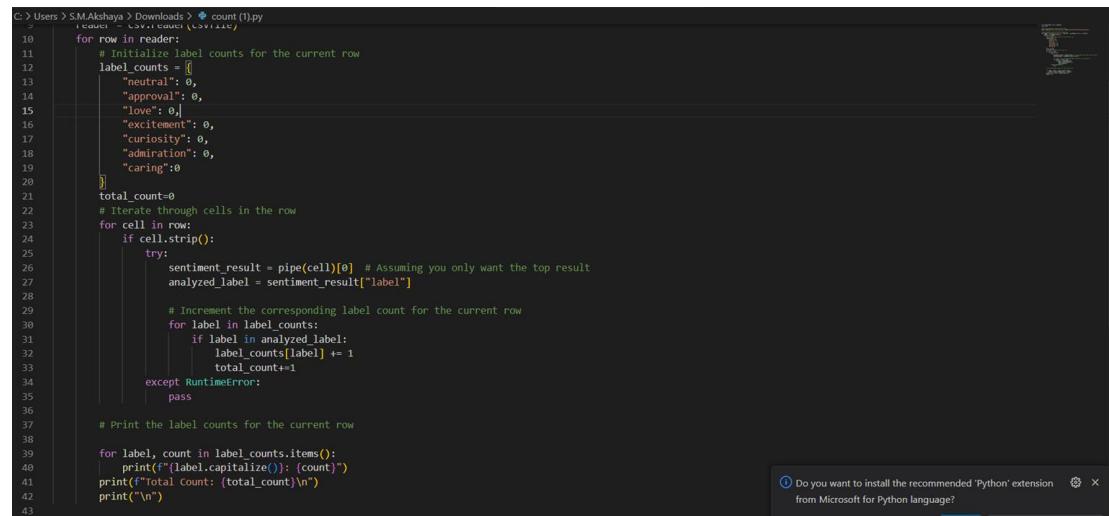
a.Sentimental Analysis Pipeline:



A screenshot of Visual Studio Code showing a Python script named 'count (1).py'. The code uses the 'transformers' library to load a sentiment classification model ('SamLowe/roberta-base-go_emotions') and iterates through rows of a CSV file ('merix1.csv') to count emotions. The interface includes a sidebar with project files like 'eee.py', 'rev (1).py', and 'from transformers import pipeline (1).py'. The bottom status bar shows tabs for PROBLEMS, OUTPUT, DEBUG CONSOLE, TERMINAL, JUPYTER, and LLDB.

```
C:\> Users > S.M.Akshaya > Downloads > count (1).py
1  from transformers import pipeline
2  import csv
3
4  # Load the sentiment classification model
5  pipe = pipeline("text-classification", model="SamLowe/roberta-base-go_emotions")
6
7  # Open the CSV file and iterate through rows
8  with open(r"D:\py\merix1.csv", "r", newline="", encoding="utf-8") as csvfile:
9      reader = csv.reader(csvfile)
10     for row in reader:
11         # Initialize label counts for the current row
```

b.Label Counts and Analysis:



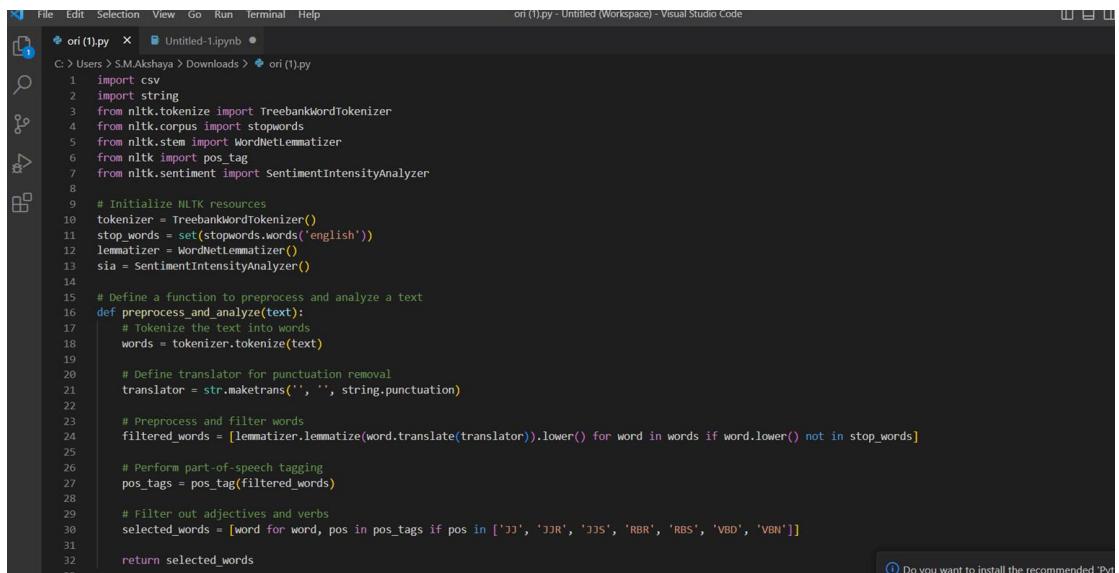
A screenshot of Visual Studio Code showing the detailed Python script 'count (1).py'. The script reads a CSV file and counts various emotions for each row. It uses a try-except block to handle errors and prints the total count of each emotion. A tooltip at the bottom right suggests installing the 'Python' extension from Microsoft.

```
C:\> Users > S.M.Akshaya > Downloads > count (1).py
10  for row in reader:
11      # Initialize label counts for the current row
12      label_counts = [
13          "neutral": 0,
14          "approval": 0,
15          "love": 0,
16          "excitement": 0,
17          "curiosity": 0,
18          "admiration": 0,
19          "caring": 0
20      ]
21      total_count=0
22      # Iterates through cells in the row
23      for cell in row:
24          if cell.strip():
25              try:
26                  sentiment_result = pipe(cell)[0] # Assuming you only want the top result
27                  analyzed_label = sentiment_result["label"]
28
29                  # Increment the corresponding label count for the current row
30                  for label in label_counts:
31                      if label in analyzed_label:
32                          label_counts[label] += 1
33                  total_count+=1
34              except RuntimeError:
35                  pass
36
37      # Print the label counts for the current row
38
39      for label, count in label_counts.items():
40          print(f"\t{label.capitalize()}: {count}")
41      print(f"\tTotal count: {total_count}\n")
42      print("\n")
```

Counting the number of happy emotions that got labelled as a whole for every product.

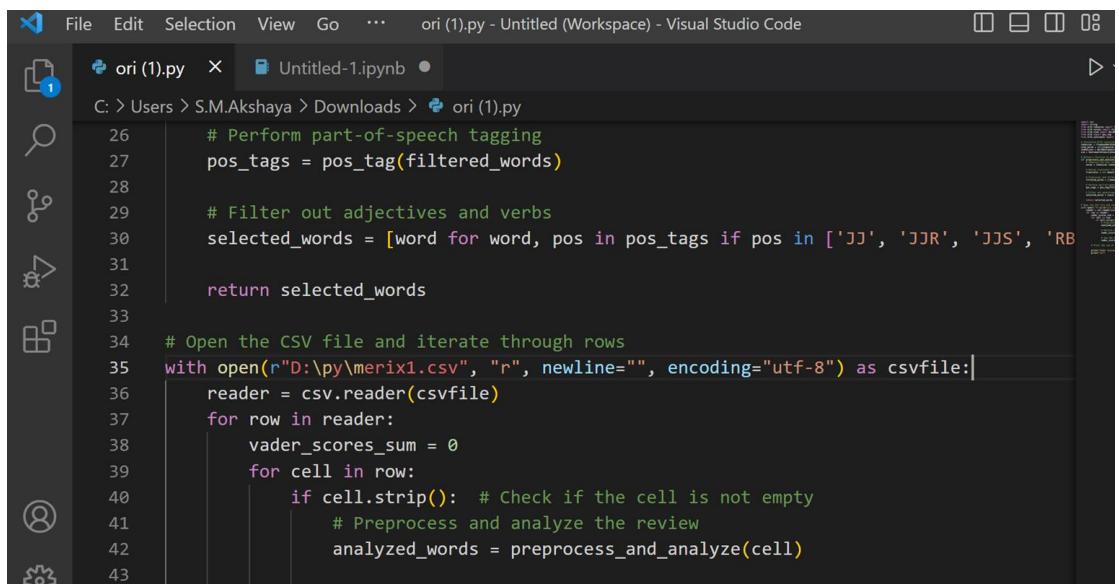
4.1.6 Summarized Sentiment Scoring of reviews

a.Data Reading and Preparation:



```
ori (1).py - Untitled-1.ipynb • ori (1).py
C: > Users > S.M.Akshaya > Downloads > ori (1).py
1 import csv
2 import string
3 from nltk.tokenize import TreebankWordTokenizer
4 from nltk.corpus import stopwords
5 from nltk.stem import WordNetLemmatizer
6 from nltk import pos_tag
7 from nltk.sentiment import SentimentIntensityAnalyzer
8
9 # Initialize NLTK resources
10 tokenizer = TreebankWordTokenizer()
11 stop_words = set(stopwords.words('english'))
12 Lemmatizer = WordNetLemmatizer()
13 sia = SentimentIntensityAnalyzer()
14
15 # Define a function to preprocess and analyze a text
16 def preprocess_and_analyze(text):
17     # Tokenize the text into words
18     words = tokenizer.tokenize(text)
19
20     # Define translator for punctuation removal
21     translator = str.maketrans('', '', string.punctuation)
22
23     # Preprocess and filter words
24     filtered_words = [lemmatizer.lemmatize(word.translate(translator)).lower() for word in words if word.lower() not in stop_words]
25
26     # Perform part-of-speech tagging
27     pos_tags = pos_tag(filtered_words)
28
29     # Filter out adjectives and verbs
30     selected_words = [word for word, pos in pos_tags if pos in ['JJ', 'JJR', 'JJ$S', 'RBR', 'RBS', 'VBD', 'VBN']]
31
32     return selected_words
33
34
35
36
37
38
39
40
41
42
43
```

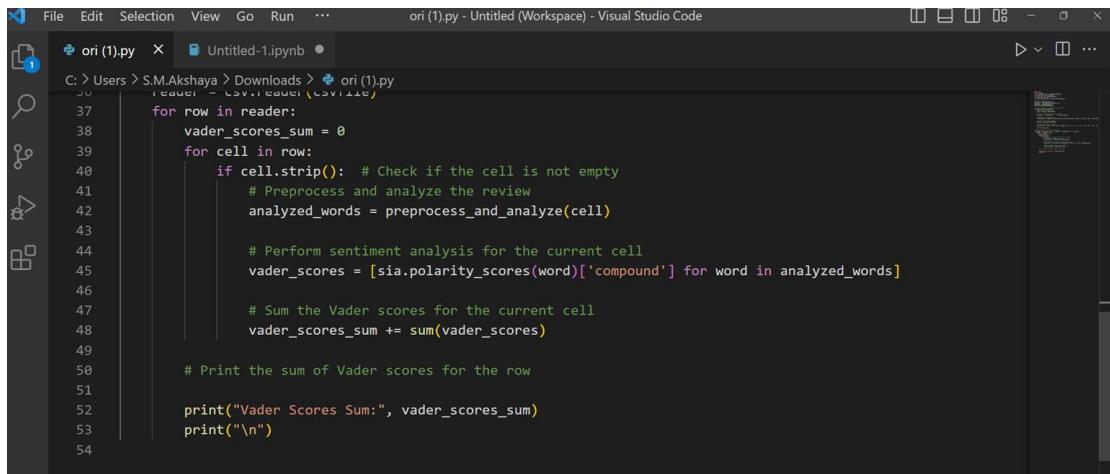
b.Text Preprocessing:



```
ori (1).py - Untitled-1.ipynb • ori (1).py
C: > Users > S.M.Akshaya > Downloads > ori (1).py
26     # Perform part-of-speech tagging
27     pos_tags = pos_tag(filtered_words)
28
29     # Filter out adjectives and verbs
30     selected_words = [word for word, pos in pos_tags if pos in ['JJ', 'JJR', 'JJ$S', 'RBR', 'RBS', 'VBD', 'VBN']]
31
32     return selected_words
33
34     # Open the CSV file and iterate through rows
35     with open(r"D:\py\merix1.csv", "r", newline="", encoding="utf-8") as csvfile:
36         reader = csv.reader(csvfile)
37         for row in reader:
38             vader_scores_sum = 0
39             for cell in row:
40                 if cell.strip(): # Check if the cell is not empty
41                     # Preprocess and analyze the review
42                     analyzed_words = preprocess_and_analyze(cell)
43
```

Using tokenizer and pos_tag which helps in filtering out the part of speech that are mentioned and we are eliminating the stopwords that are not necessary in the sentimental analysis.

c.Sentiment Analysis:

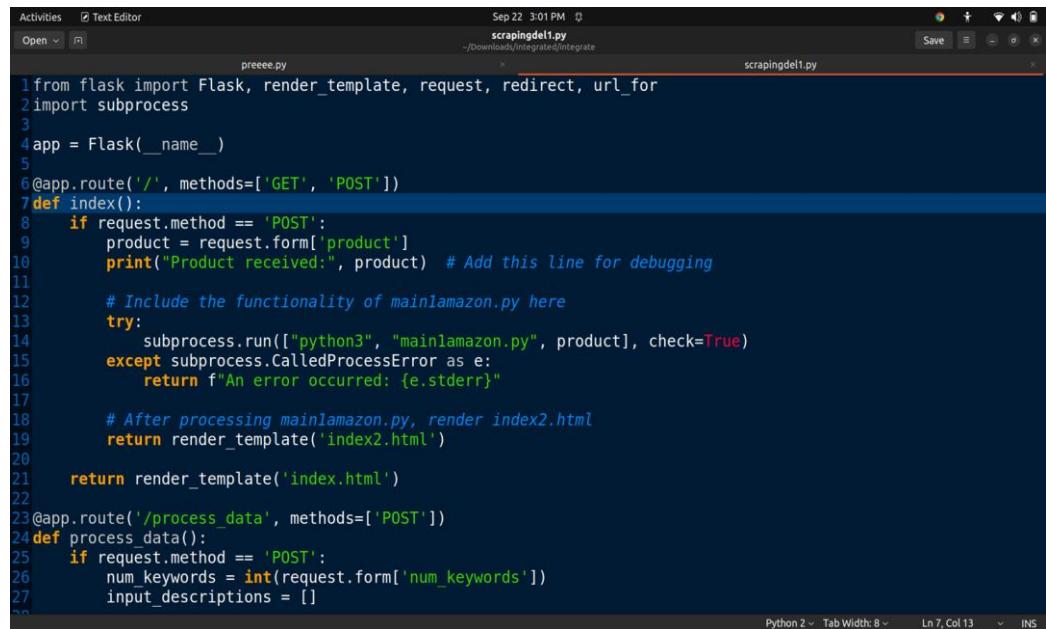


A screenshot of Visual Studio Code showing a Python script named 'ori (1).py'. The code reads a CSV file 'reader' and calculates the sum of Vader scores for each row. The code is as follows:

```
C: > Users > S.M.Akshaya > Downloads > ori (1).py
37     for row in reader:
38         vader_scores_sum = 0
39         for cell in row:
40             if cell.strip(): # Check if the cell is not empty
41                 # Preprocess and analyze the review
42                 analyzed_words = preprocess_and_analyze(cell)
43
44                 # Perform sentiment analysis for the current cell
45                 vader_scores = [sia.polarity_scores(word)['compound'] for word in analyzed_words]
46
47                 # Sum the Vader scores for the current cell
48                 vader_scores_sum += sum(vader_scores)
49
50             # Print the sum of Vader scores for the row
51             print("Vader Scores Sum:", vader_scores_sum)
52             print("\n")
53
54
```

Finally to choose , we use the vader score which makes a score out of the sentimental analysis hence the scores are made to calculate and find the sum of vader scores of each row which contains all the reviews of a particular product.

4.1.7 Frontend using Flask:



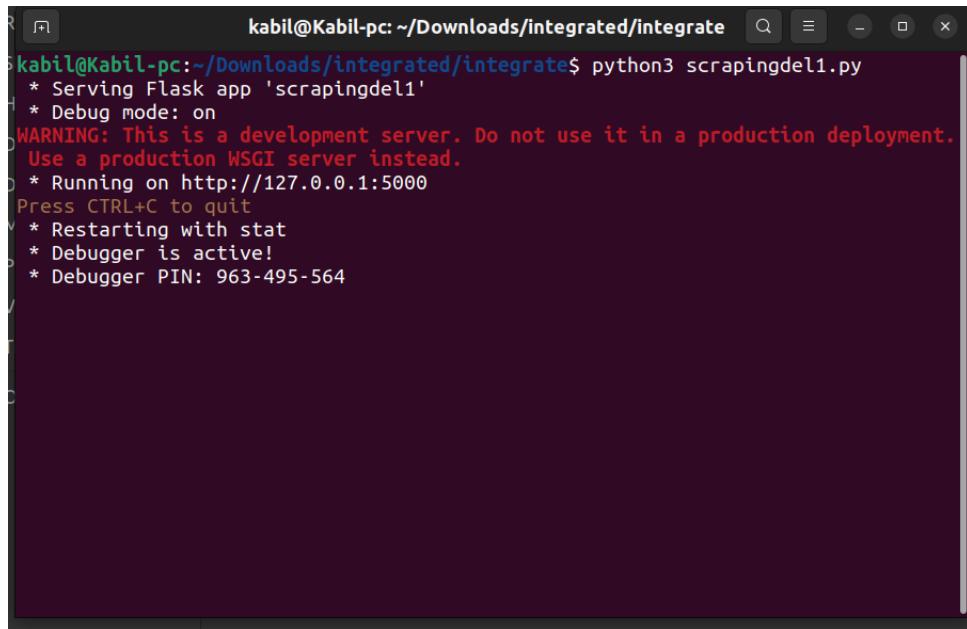
A screenshot of a terminal window showing a Python script named 'preeee.py'. The script uses Flask to handle POST requests and run a subprocess command. The code is as follows:

```
Activities Text Editor Sep 22 3:01 PM
Open preeee.py
scrapingel1.py
-/Downloads/integrated/integrate
Save
scrapingel1.py
1from flask import Flask, render_template, request, redirect, url_for
2import subprocess
3
4app = Flask(__name__)
5
6@app.route('/', methods=['GET', 'POST'])
7def index():
8    if request.method == 'POST':
9        product = request.form['product']
10       print("Product received:", product) # Add this line for debugging
11
12       # Include the functionality of mainlamazon.py here
13       try:
14           subprocess.run(["python3", "mainlamazon.py", product], check=True)
15       except subprocess.CalledProcessError as e:
16           return f"An error occurred: {e.stderr}"
17
18       # After processing mainlamazon.py, render index2.html
19       return render_template('index2.html')
20
21   return render_template('index.html')
22
23@app.route('/process_data', methods=['POST'])
24def process_data():
25    if request.method == 'POST':
26        num_keywords = int(request.form['num_keywords'])
27        input_descriptions = []
28
```

Finally to integrate all these backend with frontend website,Flask has been used.

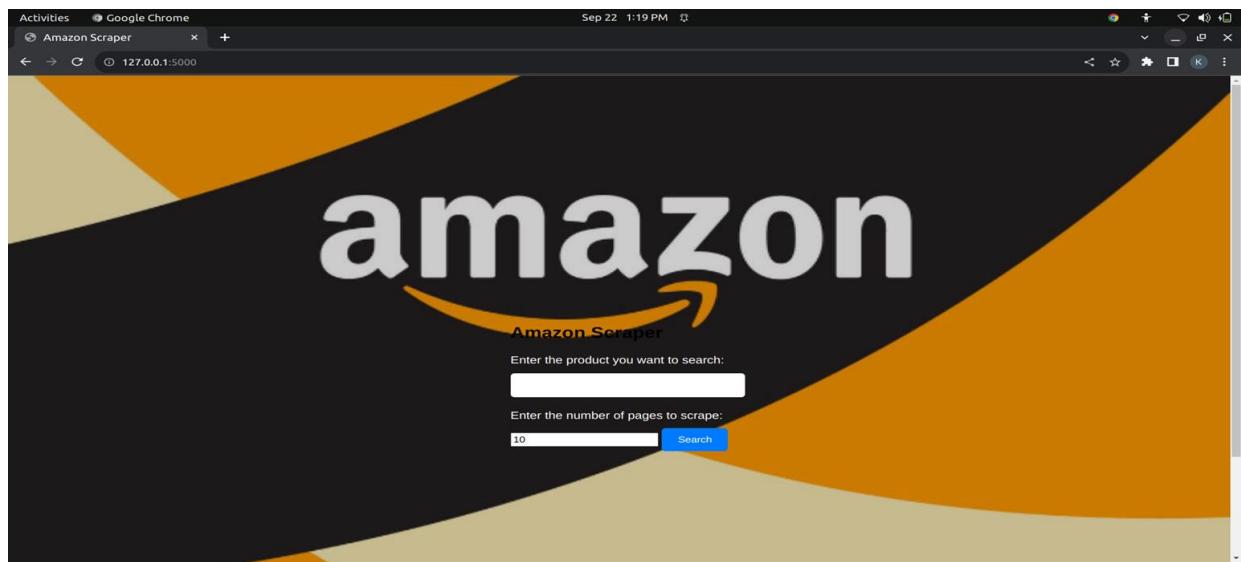
4.2 RESULTS

4.2.1 Web Scraping:-

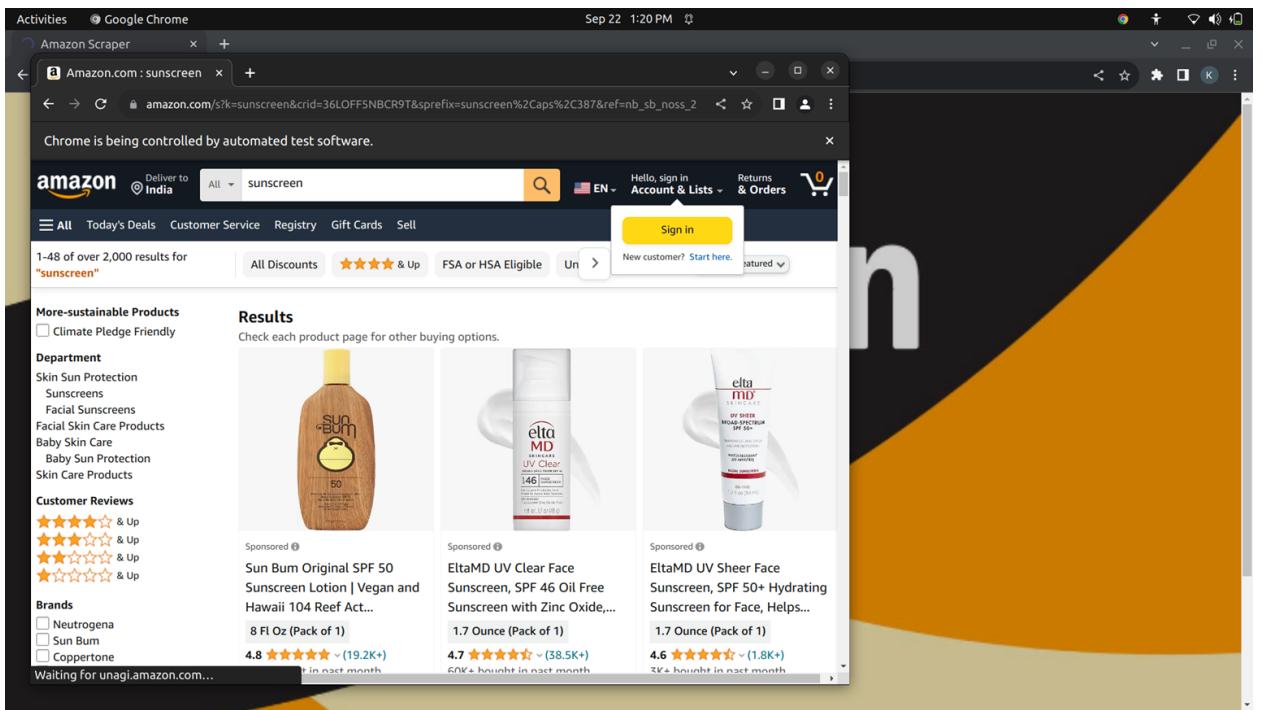


A terminal window titled 'kabil@Kabil-pc: ~/Downloads/integrated/integrate'. It shows the command 'python3 scrapingdel1.py' being run, followed by the output of the Flask development server starting up. The output includes: '* Serving Flask app 'scrapingdel1'' , '* Debug mode: on' , 'WARNING: This is a development server. Do not use it in a production deployment.' , 'Use a production WSGI server instead.' , '* Running on http://127.0.0.1:5000' , 'Press CTRL+C to quit' , '* Restarting with stat' , '* Debugger is active!' , '* Debugger PIN: 963-495-564'.

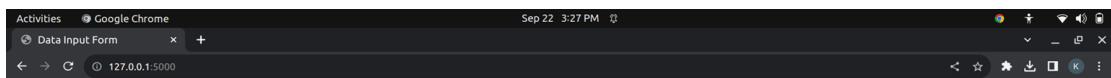
We have used the linux environment to the run the program, thus starting with the scrapping the data from the amazon site.



The front-end that we have incorporated with the flask which deploys where the user gives the input of the product name which is fed in input of amazon's search box.



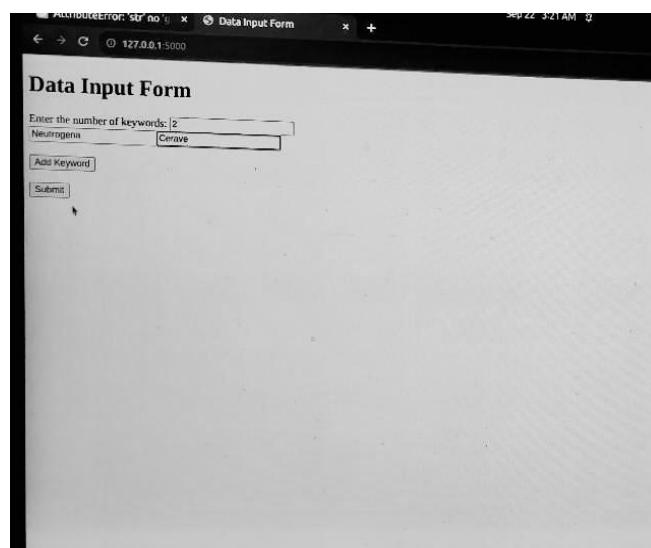
This shows the scrapping of data which almost scraps nearly 1200 rows of data approximately, fetching the product's description, url, overall rating , reviewcount and the overall rating of the products , this almost scraps the datas of 21 pages from the Amazon's site.



Data Input Form

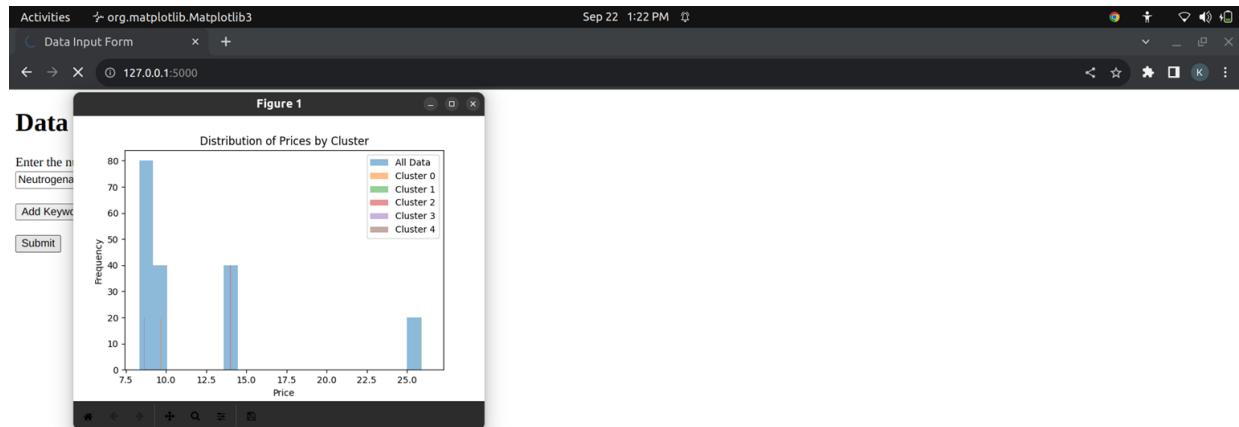
Enter the number of keywords:

This page asks the user for the product's specifications as inputs by giving the number of inputs as a first step, where the code filters and saves the data rows accordingly to the given specifications.
And hence we use the saved present file for the upcoming processes.



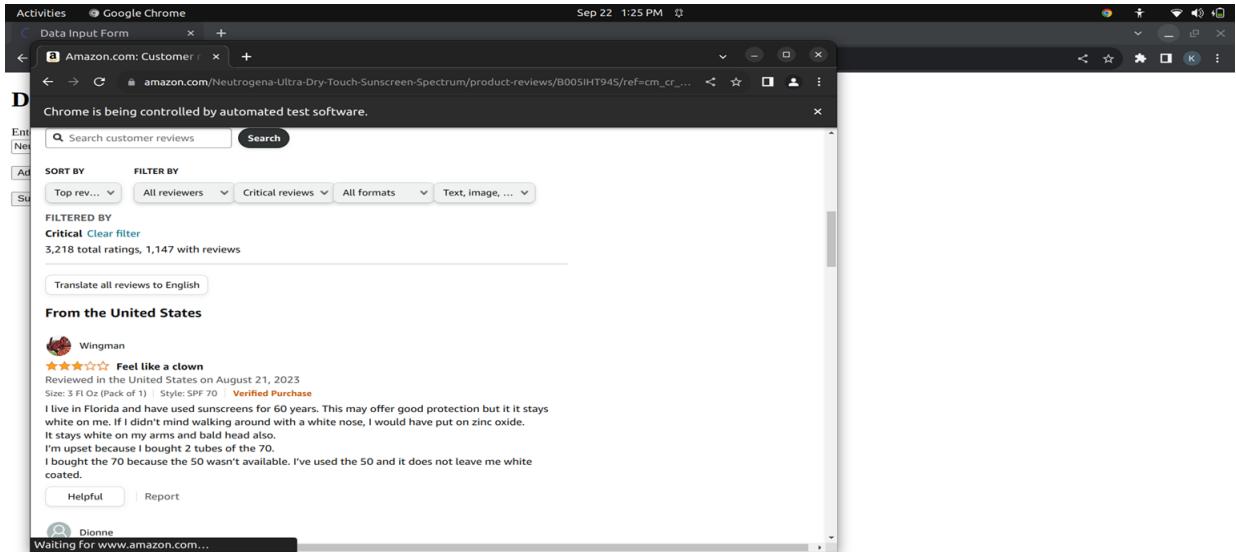
User specifying the inputs

4.2.3 Distribution of Prices:-



This picture shows the distributions of prices visually in graphs and which are classified by k-means clustering algorithm considering products with high review counts and grouping them into ranges and considering the price ranges where most of the products lie in ,and filtering the datas in a separate csv and proceeding with the next process for filtering out.

4.2.4 Review Scrapping:-



Scraping all the reviews , even the critical reviews which are already categorized by Amazon and saving them in a csv , we scrap almost 40-50 reviews for each product and classifying them into emotion using the Roberta pretrained model which classifies them into emotions : admiration, amusement, anger, annoyance, approval, caring, desire, confusion, curiosity, disgust, excitement, disappointment, joy, love and more on.

And counting the happy emotions alone i.e., neutral, approval, love, excitement, curiosity, admiration, caring for each product and classifying them which gets the higher score.

And finally taking the reviews of those filtered and removing all the stop words and using lemmatizers and tokenizers from the “Natural Language Tool kit” and even performing part of speech tagging with the reviews i.e., taking /considering all the adjectives and verbs and calculating the vader sum and hence by comparing the resultant vader sum with all the products , we get the best recommended product that many people were happy about it and considering using it.

4.2.5 Final output:-

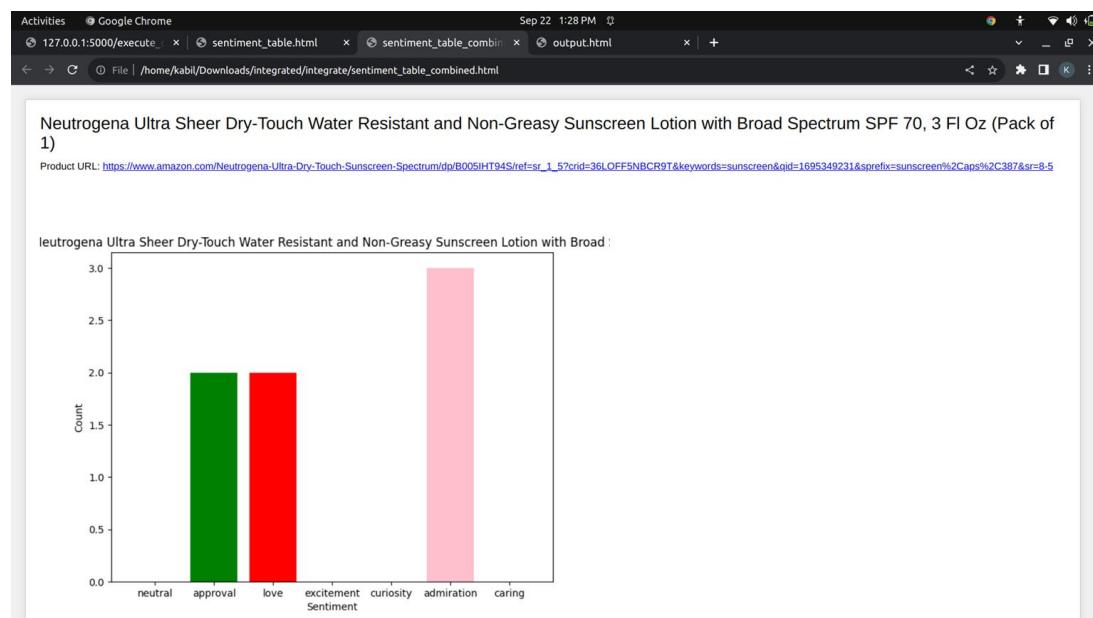
Activities Google Chrome
Data Input Form sentiment_table.html +
File /home/kabli/Downloads/integrated/integrate/sentiment_table.html

Neutrogena Ultra Sheer Dry-Touch Water Resistant and Non-Greasy Sunscreen Lotion with Broad Spectrum SPF 70, 3 Fl Oz (Pack of 1)
https://www.amazon.com/Neutrogena-Ultra-Dry-Touch-Sunscreen-Spectrum/dp/B005IHT94S/ref=sr_1_5?crid=36LOFF5NBR9T&keywords=sunscreen&qid=1695349231&sprefix=sunscreen%2Caps%2C387&sr=8-5

Sentiment	Score
Error	
[{"label": "love", "score": 0.4864284098148346}]	
[{"label": "admiration", "score": 0.6772761940956116}]	
[{"label": "approval", "score": 0.5449449419975281}]	
[{"label": "admiration", "score": 0.6318193674087524}]	
[{"label": "joy", "score": 0.3173394501209259}]	
[{"label": "admiration", "score": 0.5581114888191223}]	
[{"label": "approval", "score": 0.4692990481853485}]	
[{"label": "love", "score": 0.8241609930992126}]	
[{"label": "joy", "score": 0.6096246242523193}]	

Neutrogena Beach Defense Water-Resistant Face & Body SPF 70 Sunscreen Lotion with Broad Spectrum UVA/UVR Protection Oil-Free Fast-Absorbing Sunscreen Lotion

This shows the labelled classification of emotions using Roberta Model.



This shows that admiration got labelled on more number of reviews.

Product Vader Scores		
	Product URL	Vader Scores Sum
Neutrogena Ultra Sheer Dry-Touch Water Resistant and Non-Greasy Sunscreen Lotion with Broad Spectrum SPF 70, 3 Fl Oz (Pack of 1)	Link	9.8581
Neutrogena Beach Defense Water-Resistant Face & Body SPF 70 Sunscreen Lotion with Broad Spectrum UVA/UVB Protection, Oil-Free Fast-Absorbing Sunscreen Lotion, Oxybenzone-Free, 6.7 oz	Link	2.9180
Neutrogena Sport Face Sunscreen SPF 70+, Oil-Free Facial Sunscreen Lotion with Broad Spectrum UVA/UVB Sun Protection, Sweat-Resistant & Water-Resistant, 2.5 fl. oz	Link	11.8622
Neutrogena Ultra Sheer Dry-Touch Sunscreen Lotion, Broad Spectrum SPF 55 UVA/UVB Protection, Lightweight Water Resistant Face & Body Sunscreen, Non-Greasy, Travel Size, 3 fl. oz	Link	0.5592

The shown vader scores are sum of the review's it is , by which we can conclude that the higher vader score id ,being great the product is with positive appreciation given by users.

5.1 CONCLUSION:

The usage of webscrapping by the automation of browser which in order collects the data at that instance available in the e-commerce site , and the data can be fetched and analysed at any time with the usage of webscrapper language comprising Selenium and Beautiful Soup with python.

After several filtrations of data containing the details of product and reviews , we are able to get the best products finalised list based on the ranking and vader scores .

Hence ,this resulting in an application which can be used to find the best product url list according to the specifications given as the input and produces results with vaderscores (higher the vader scores are implies that higher the good product it is with more good reviews).

This application saves the user a lot of time and one does not want spend hours on going through the list of product available on the e-commerce sites, searching for the desired they need with good reviews and meeting within their budget.

Here are some results given by our model:-

- 1)Sunscreen: 4.96 (Cerave)
- 2)Laptop: 4.24 (Lenovo)

5.2 FUTURE WORK:

Collecting product data from Amazon through web scraping and then determining the best product URL based on ratings and sentiment analysis of reviews has several potential applications and future prospects:

- 1) Competitor Analysis: You can track competitors' product ratings and reviews over time to gain insights into market trends and how they are perceived by customers.
- 2) Product Quality Assessment: Manufacturers and sellers can use this data to assess the quality of their products by analyzing customer feedback and making necessary improvements.
- 3) Customized Shopping Experiences: Leverage the data to offer personalized shopping experiences, suggesting products based on individual preferences and reviews they've given in the past.
- 4) Sentiment-Driven Marketing Campaigns: Develop marketing strategies based on sentiment analysis. For example, launch targeted advertising campaigns for products with overwhelmingly positive reviews and ratings.
- 5) Trend Forecasting: Analyze Amazon product reviews to identify emerging trends and consumer preferences. This information can be valuable for retailers and manufacturers planning their product lines.
- 6) Price Comparison Websites: You can use this data to build or enhance price comparison websites. Users often seek the best deals, and your analysis can help them find not only the cheapest products but also those with high ratings and positive sentiment in reviews.
- 7) Market Research: Analyzing sentiment from product reviews can provide insights into customer opinions and preferences. Companies can use this data to improve their products, marketing strategies, and customer service.

5.3References:

1. “*EFFICIENT SCRAPING OF DATA FROM WEBSITES USING SELENIUM*”

Authors: [1]Shreya V. Dhone student , [2] Anupama D. Sakhare Asst Prof,[3] Satish J.Sharma Prof.

1,2,3 Dept of Electronics and Computer Science, Rakshtrasant Tukdoji Maharaj Nagpur University, India . 2022 jetir, volume 9, issue 6, issn:2349-5162

<https://www.jetir.org/papers/JETIRFM06063.pdf>

2. “*RoBERTa-GRU: A Hybrid Deep Learning Model for Enhanced Sentiment Analysis*”

Authors: Kian Long Tan , Chin Poo Lee, Kian Ming Lim , Multimedia University, Malaysia

<https://doi.org/10.3390/app13063915>

3. “*The k-means Algorithm: A Comprehensive Survey and Performance Evaluation*”

Authors: Mohiuddin Ahmed, Raihan Seraj, Syed Mohammed Shamsul Islam

<https://www.mdpi.com/2079-9292/9/8/1295>

4. “*Sentiment analysis of product reviews: A review*”

Authors: T.K. Shivaprasad, Jyothi Shetty, NMAM Institute of Technology

<https://ieeexplore.ieee.org/document/7975207>

5. “*A review on sentiment analysis and emotion detection from text*”

Authors: Nandwani, Rupali Verma

<https://link.springer.com/article/10.1007/s13278-021-00776-6#article-info>

6. “*VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text*”

Authors: C.J. Hutto, Eric Gilbert , Georgia Institute Of Technology

<http://eegilbert.org/papers/icwsm14.vader.hutto.pdf>

7. “*An overview and comparison of free Python libraries for data mining and big data analysis*”

Authors: Stacin, Jovic

<https://ieeexplore.ieee.org/document/8757088>

8. “*An Automated Web Application Testing System*”

Authors: Tarek M Mahmoud, Moheb Girgis, Bahgat A. Abudullatif, Alaa Zaki

https://www.researchgate.net/publication/270569315_An_Automated_Web_Application_Testing_System