Hardik Prabhu
241058019

1> Describe the Primary Purpose of User Aget header in web Scrappis.

→ *The User Agent header is an HTTP header used by a client (browser or bot) to identify itself to the server. When a server Recieve a roqust it can check the User-Aget String to determine what kind of device or software is makes the Roqut This helps to serve tailor the contet to differt devics or browers.

* In web srappr setting a User-Agt headr is essetid to mimic a real browser, makes your sroper look less likes to be blocked.

* So the are the primay pupoes of user agt header,

2> Wolt s the purpse of robots.txt file on a website.

→ The robots.txt file is a text file placed in the root directoy of a website that provids guidlne for web crauh and bot about which pats of the site they are allowed to accen. It uses the robot Exchion Protcol & can include rules to restrict specific crawbs or path on the site.

Ex:-
    User-aget: *

Disallow : /admin /
Disallow : /private -data /

In the above example , all bod (User-agent : *) are
instructed not to acces the /admin ( and /private -data/
directors. While robots.txt is not a security factors

3) Scrape data contents from single web page using
Beautiful Soup library.

→ Here the website used is http://quotes.toscrape.com/

```
from bs4 import Beautiful Soup
import requests

url = "http://quots.tosrape.com/"
response = requests.get(url)
soup = Beautiful Soup (response.content, "html.parser')
quots = soup.find_all ("div", class_ ="quots")
for quote in quotes:
    text = quote.find ("span", class_="text" ).get_text()
    author = quote.find ("small", class_="author" ).get_text()
    print (f"Quote : {text} \Author : {author}\n")
```

4) Scrape data contents navised throgh multiple web
pages using Beautiful Soup library.

→
```
from bs4 import Beautiful Soup
import Requests
base_url = 'http://quots.tosrape.com/page/'
```

```
for page in range (1,6):
    url = f"{base_url}{page}/"
    response = requests.get(url)
    soup = BeautifulSoup(response.content, "html.parser")
    quotes = soup.find_all("div", class="quote")
    for quote in quotes:
        text = quote.find("span", class="text").get_text()
        author = quote.find("small", class="author").get_text()
        print(f"Page {page} - Quote : {text}\nAuthor : {author}\n")
```

5) Describe the Architecture of Scrappy :-

→ The Architecture of Scrapy is the follow compots

   * Scrapy Engine :-

   The Engine is responsible for controlls the
data flow between all compots of the system. and
trigger evets when contain actions occur

   * Scheduler :-

   The Schedule recievs requests from the engine and
enqueu them for proody them later. when
the Engine Request them.

   * Downloader

   The Downloads is Responsible for fetch web
pages and fooch them to the Engine which in turn
feeds them to the Spider.

* **Item pipelines**

The item pipeline is responsible for proce
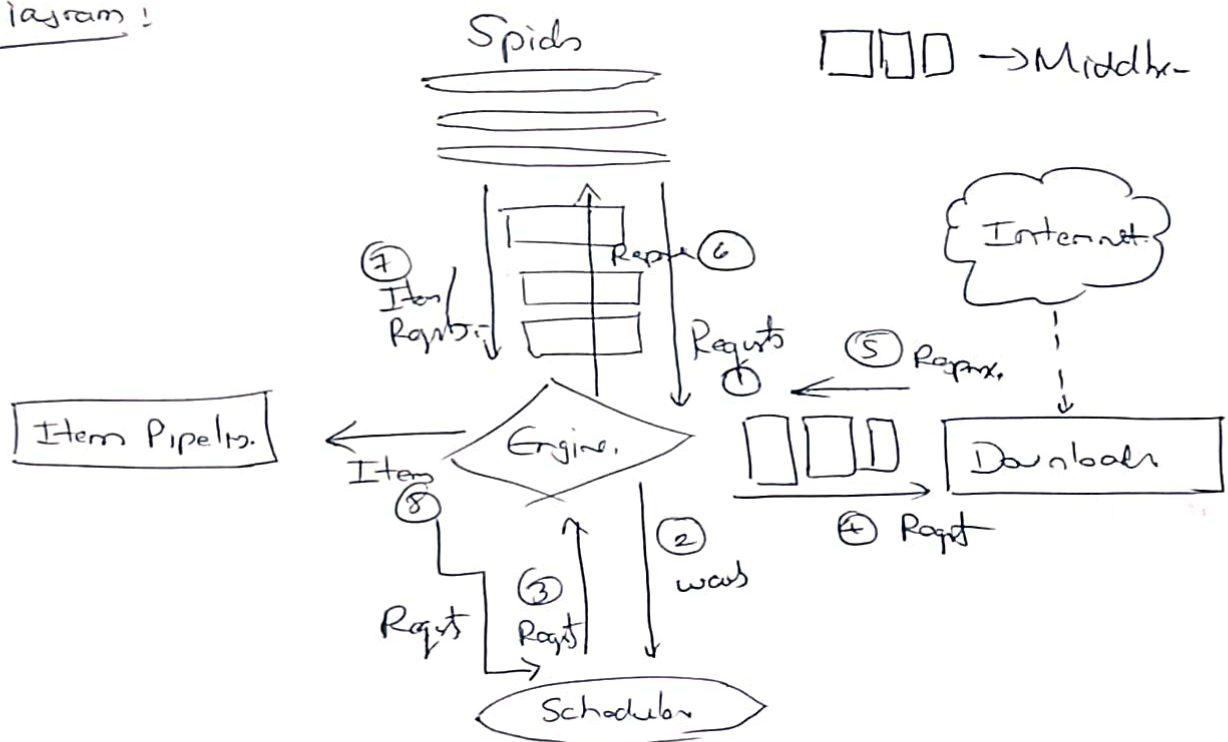the items once they have been extracted by
the spiders.

* **middlewares :-**

• **Downloader middleware** :- These are the specific hooks
that are sitting between the engine & Downloader.

• **Spider Middleware :-**

These are the specific hooks that sit between
Engine & the Spiders and are able to process
Spider Input and output.

**Diagram :**



So this is the Architect of the Scrapy
framework used for web scraping.

6) Demonstrate Scrapy framework for the proxy
Example:-

```
import Srapy
class ProductSpider (srapy.spider):
    name : "product_spider"
    start-urls > ["http://www.spom-Tr.com]
    def parse (self, response):
        for product in response.CSS ('product"):
            yield {
                "name": product.css ('name ::text ').get()
                "cost": product.css ('cost ::text ").get()
                "manufach": product.css ('manufach ::text").get()
                "ratng": product.css ('ratn ::text').get()
            }
```

        Now we need to run th file
    by, srapy runspider s1 product_spider.py -o
                                products.json.

7) Creats a pdas dataframe from Super market
    Invoice Data, & exects.

→   i) Loads

        import pars s pd.
        df = pd. cread_csv ('Supermarket-invoice_data.csv')

        a) Cont the number of customs per product
            category or depatmt

customer_count = df.groupby("Description")["Customer Code"].nunique())

print(customer_count)

b) Calculate the average customer visit per product category

avg-visit = df.groupby("Description")["Quantity"].mean()

c) Identify the product brought by the most customers

popular_product = df.groupby("Description")["Customer Code"].nunique.idmax()

print(popular_product)

d) Sort Customer based on age, product Purchased, or product cart.

Sorted_cust = df.sort_values(by = ["Sales", "Quantity"], ascending = [False, False])

print(Sorted_customer)

i) Explain Category & Explanatory variations based on the Relationship between three necessary Players.

* Univariates
  * It shows the distribution of Single variable
  * Visual example
    • Histogram
    • Bar plots.

* Bivariate:-
  * Shows the Relationship between two variable
    • Visual Example

- Scatter plots, box plots.

× Multivariate:
  - Show relationships among multiple variables.
  - $\underline{vib}$ (e.g heatmaps, pair plots)

2) Define DAX, Demonstrate 10 DAX functions with context.

→ DAX which could be called s Data Analysis Expression is a formula language for data calculation in Power BI, SSAS & Excel Power Pivot

- Sum :
  - TotalSales = SUM (Sal [Amount])
- Average.
  - Avg = AVERAGE (Sales [Amount])
- Count
  - Count = Count (Sales [CustomerID])
- Count Rows
  - CountRows = Total Rows = COUNTROWS (Sales)
- Min
  - Min Sales = MIN (Sales [Amount])
- Max
  - Max Sales = MAX (Sales [Amount])
- RELATED
  - Related Category = RELATED (Product [category])
- IF
  - Sales Status = IF (Sales [Amount] > 500, "High", "Low")
- Calculate
  - High Sales = CALCULATE (SUM (Sales [Amt]), Sales [Amount] > 100)

- FILTER

Filtered = FILTER (Sal , Sal [Amt ]>500)

3) Differ between Calculated Measure & Calculated Column.

| Calculated Measure | Calculated Column |
|---|---|
| * Used for Aggregation & Dynamic Calculations that chu based on fib & report Context | * Used to create additional columns with static value in a table, calch row by row |
| * Computed dynamically based on the report filter cont | * Calc once during data fresh, vals do not chu with filt |
| * Stored as a formula in the model & Calculles on demand | * Stored as actual value in the data model. |
| * Cannot be placed directly | * Can be added directly |
| * Calculated at run time | * Pre calculated. |
| * Typically Faster | * Remains Constant until the next refresh. |
| * Calculated based on the Entire Dataset. | * Limited to row level Context |
| * Uses for total Sal aver matrics etc. | * Used for creation category and more |
| * Uses functions like SUM() AVERAGE() etc. | * Often use IF(), CONCATENATE() etc |
| * Can be used in Measure, cards KPI's etc | * Can be used in table, matrices etc. |