

# PDV Assignment

Name: Karthik  
Reg No: 24 1058016  
Big Data Analytics

1) Describe the primary purpose of a User-Agent header in web scraping.

→ User-Agent header in web scraping plays a crucial role in mimicking the behavior of a real web browser.

Primary purpose of the User-Agent Header:

1. Simulate a Real User/Browser

The User-Agent header tells the server what type of device or browser is making the request. When scraping, you can set this header to resemble a real browser (like Chrome, Firefox etc) to prevent the server from blocking the request as suspicious or coming from an automated tool.

2. Bypass Anti-Scraping Mechanisms

• Many websites use anti-bot measures that detect and block requests with non-standard or missing User-Agent headers, as these are often associated with web scrapers or bots.

3. Handle Device or Browser-Specific Content

→ Some websites deliver different content depending on the browser or device requesting it. By including a User-Agent header that mimics a specific browser.

4. Respect Website's Terms and Conditions

→ While not a guarantee of respect, using a valid User-Agent header that identifies your scraper as a specific browser can indicate to the server that the request is coming legitimate.

2) what is the purpose of a robots.txt file on a website?

→ The robots.txt file is a text file placed in the root directory of a website that provides guideline for web crawlers' access to website. It uses the robots exclusion protocol & can include rules to restrict specific crawl length or path on the site.

Ex:- user-agent: \*

Disallow: /admin/

Disallow: /private/

\*Disallow :- Indicates which pages or directories should not be crawled by the specified user-agent. If the value is left empty (i.e. Disallow), it means the crawler is allowed to access all parts of the site.

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

→ Here's an example of how to scrape data from a webpage

```
from bs4 import BeautifulSoup  
import requests
```

```
url = "http://quotes.toscrape.com/"
```

```
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()
```

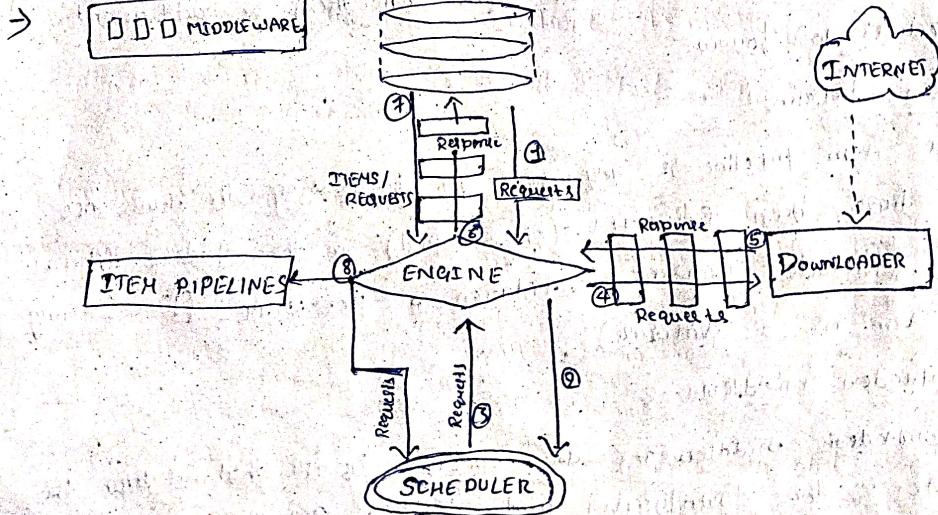
```
    author = quote.find("small", class_="author").get_text()
```

```
    print(f"Quote : {text} | Author : {author} )"
```

4. Scrape data contents navigating through multiple web page using BeautifulSoup library.

```
→ from bs4 import BeautifulSoup  
import requests  
base_url = "http://quotes.toscrape.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 Scrapy framework.



The data flow in Scrapy is controlled by the Execution Engine, and goes like this:

#### \* Scrapy Engine:-

The Engine is responsible for controlling the data flow b/w all components of the system, and triggering events when certain action occurs. See the Data flow section above for more detail.

#### \* Scheduler:-

The scheduler receives requests from the Engine, and enqueues them for feeding them later (also to the Engine) when the Engine requests them.

#### \* Downloader:-

The Downloader is responsible for fetching web pages and feeding them ~~to~~ to the Engine which, in turn, feeds them to the Spider.

#### \* Spider

Spiders are custom classes written by Scrapy users to parse responses & extract items from them or additional requests to follow.

#### \* Item pipeline

The item pipeline is responsible for processing the items once they have been extracted (or scraped) by the Spiders. Typical tasks include cleaning, validation & persistence (like storing the item in a database).

#### Downloader middlewares

Downloader middlewares are specific hooks that sit b/w the Engine & the downloader and processes requests when they pass from the Engine to the Downloader, & responses that pass from Downloader to the Engine.

## Spider middleware

Spider middlewares are specific hooks that sit b/w the Engine & the Spiders and are able to process spider input (responses) and output (items and requests).

6. Demonstrate Scrapy framework with an example to store information about a product having attributes.

→ Set up Scrapy project

Scrapy startproject shopping-Scraper

It will create a folder structure

Inside the shopping-Scraper/items.py.

Import Scrapy

class ProductItem (Scrapy.Item)

name = Scrapy.Field()

cost = Scrapy.Field()

manufacturer = Scrapy.Field()

rating = Scrapy.Field()

shopping-Scraper/spider → create new python file. Product-Spider.py.

import Scrapy

from shopping\_Scraper.items import ProductItem

• class ProductSpider (Scrapy.Spider):

name = 'Product-Spider'

start\_urls = ['https://www.shopping-tv.com']

def parse(self, response):

for product in response.css('product'):  
yield {

"name": product.css('name ::text').get()

"cost": product.css('cost ::text').get()

"manufacturer": product.css('manufacturer ::text').get()

```
"rating": product.css("rating ::text").get()
```

```
}
```

```
next_page = response.css('a.next-page::attr(href)').get()
```

```
if next_page:
```

```
    yield response.follow(next_page, self.parse)
```

7) Create a pandas dataframe from Super Market Invoice Data, and execute following operations.

i) Load the data

```
import pandas as pd
```

```
df = pd.read_csv("Supermarket-invoice-data.csv")
```

a) Count the number of customers per product category or department.

```
customers_count = df.groupby('Category') ['customer ID'].nunique()
```

```
print(customers_count)
```

b) Calculate the average customer visit per product category.

```
avg_visits = df.groupby('Category') ['customer ID'].value_counts().mean()
```

```
print(avg_visits)
```

c) Identify the product bought by the most customers

```
product_bought = df.groupby('Category') ['customer ID'].nunique().idxmax()
```

```
print(product_bought)
```

d) Sort customers based on age, products purchased, or products cost.

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

```
print(sorted_customers)
```

Explain categories of Exploratory visualizations based on the relationships b/w the three necessary players.

## 1. Univariate visualizations (single Variable)

These visualizations explore the characteristics of a single variable. The relationship here is b/w a single variable & its visual representation, where the user (observer) interprets patterns.

Example:

- Histogram : Display the distribution of a continuous variable
- Bar chart : Show frequency counts or categories for a discrete variable

## 2. Bivariate visualizations (Two Variables)

These visualizations focus on the relationship b/w two variables, either numerical or categorical. The key relationship here is between two variables & their visual mapping.

## 3. Multivariate visualizations (Multiple Variables)

When exploring the relationships b/w three or more variables, the visualization often becomes more complex. The user is tasked with interpreting multiple data dimensions simultaneously. Here, the visual representation could employ color, shape, size or even animation to represent more than

Deline DAX. Demonstrate any 10 DAX formulas with a context.

1. Count :- Counts the number of non-blank rows in a column. Typically used to count the number of rows with valid entries in a dataset.

## 2. DISTINCTCOUNT:-

Counts the distinct (unique) values in a column.

Useful for identifying the number of unique entries, such as distinct customers, products, or transactions.

## 3. AVERAGE:-

Calculates the average of a column.

Used to find the mean value, such as average sales or average quantity.

## 4. MIN:-

Returns the smallest value from a column.

Useful for identifying the minimum numerical or date value, such as the lowest sales amount or earliest date.

## 5. MAX:-

Returns the largest value from a column.

Helps identify the maximum value in a dataset, such as highest sales or latest date.

## 6. SUMMARIZE:-

Returns a summary table for the requested columns and aggregated data.

Used to group data and apply aggregations like SUM or COUNT.

## 7. CALCULATE:-

Modifies the context of a calculation with filters used to calculate results under specific conditions, such as sales for a specific product category.

## 8. IF

Returns one value if a condition is true, & another if it's false. Enables conditional calculations, such as classifying sales high/low.

## 9. SUM!

Sums the values in a column

Used to add up numerical values such as sales amounts, quantities, or profits.

## 10. ~~IFERROR~~ = SUMX?

Returns a specified result if an expression for each row summing the results. Enables summation of complex expressions across rows, such as calculating total profits for each transaction.

## 11. SWITCH!

Evaluates an expression against a list of values and returns the corresponding result. Provides a simplified way to test multiple conditions, like a case statement.

Discuss ten differences b/w calculated measure & calculated column.

Calculated Measure	Calculated Column
* Calculated on the fly based on context	* Row-wise calculation stored in the model.
* Based on filters and slices context	* Calculated for each row, independent of filters.
* Not stored: calculated dynamically	* Stored in the data model.
* Used in aggregations and summaries	* Used for row-level calculations or new dimensions
* Performance Impact recalculated each time it's used	* Calculated once, stored at load time.
* Filter we cannot be used directly in filters	* Can be used directly in filters & slices.

- \* No impact on model structure
  - \* get vs Dynamic and changed with context
  - \* get operated on aggregated or summaries
  - \* use function like sum(), Average(), etc.
- \* Add new column to the model
  - \* static once calculated.
  - \* operated on individual group.
  - \* often use IF(), concatenation(), etc.