

• PDV Assignment

i) 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 behaviour of a real web browser.

Primary purpose of the user agent header

i) 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.

ii) Bypass Anti-scraping Mechanisms :

many websites use anti-bot measures that detect and block requests with non standard or missing user agent headers or there are often associated with web scrapers or bots.

iii) Handles Different 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

iv) 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 from a legitimate

Q2) 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 crawler 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.

Q3) Scrape data contents from a single web page using BeautifulSoup.

→ Here's an example of how to scrape data from a web page.

from bs4 import BeautifulSoup

import requests

url = "https://quotes.toscrape.com/"

response = request.get(url)

Soup = BeautifulSoup(response, parse="html.parser")

quote = Soup.find_all("div", class="quote")

for quote in quote:

text = quote.find("span", class="text").gettext()

author = quote.find("small", class="author").get

text()

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

Q4) Scrape data contents navigating through multiple web page using BeautifulSoup library.

→ from bs4 import BeautifulSoup

import requests

base_url = "http://quotes.toscrape.com/page/1/"

for page in range(1, 5):

url = f"{base_url}page/{page}"

response = request.get(url)

Soup = BeautifulSoup(response, parse="html.parser")

quote = Soup.find_all("div", class="quote")

for quote in quote:

text = quote.find("span", class="text").gettext()

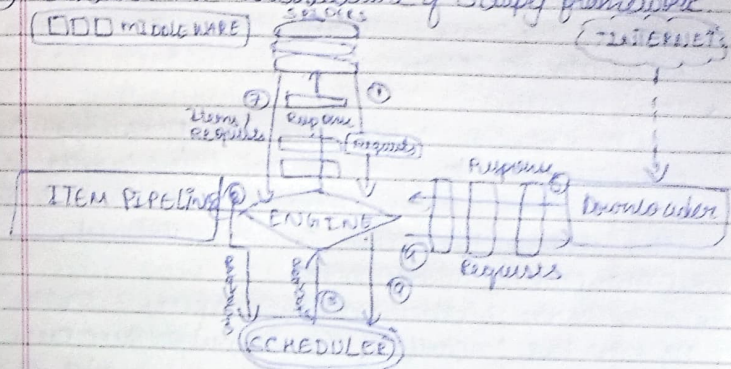
author = quote.find("small", class="author").get

text()

print(f"page {page} quote: {text} | author: {author}\n")

Q5) Describe the architecture of Scrapy framework.

□□□ MIDDLEWARE



The data flow is scrapy is controlled by the Extension Engine, and goes like this:

* Scrapy Engine:

This engine is responsible for controlling the data flow b/w an components of the system and triggering events when certain action occur. see the data flow section above for more details.

* 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 the engine which, in turn, feeds them to the spiders.

* Spiders

Spiders are custom classes written by scrapy users to parse response & extract items from them as additional requests to follow.

* Item Pipeline:

The item pipeline is responsible for processing the items and 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 process requests when they pass from the engine to the downloader & responses that pass from downloader to the engine.

Date: / /
Page No.:

* Spider middlewares

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

Date: / /
Page No.:

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/item.py.

import scrapy

class ProductItem(scrapy.Item)

name = scrapy.Field()

cost = scrapy.Field()

manufacturer = scrapy.Field()

rating = scrapy.Field()

shopping-scraper/spiders → write 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.shoppingtv.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()

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

next_page = response.css("a.next-page::text").get()

if next_page:

yield response.follow(next_page, self.parse)

2) Create a pandas dataframe from supermarket invoice data, and execute following operations

i) load the data

import pandas as pd

df = pd.read_csv("supermarket_invoice_data")

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

customer-count = df.groupby('description')

['customer_id'].unique()

print(customer-count)

b) calculate the average customer visit per product category.

avg-visits = df.groupby('description')['customer_id'].value_counts().mean()

print(avg-visits)

c) identify the products bought by the most customers

product-bought = df.groupby('description')

['customer_id'].unique().count()

print(product-bought)

d) start customer based on age, products purchased or products cost

sorted-customers = df.sort_values(by=['age', 'products', 'cost'])

ascends = (False, False)

print(sorted-customers)

3) Explain categories of exploratory visualization based on the relationships b/w the three necessary players.

i) 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 which the user (observer) interprets patterns.

examples:

Histogram: Display the distribution of a continuous variable.

Bar: Show frequency counts or categories for a discrete variable.

ii) bivariate visualizations (Two Variables)

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

iii) 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. Table Page No. Representation could employ color, shape size or even animation to represent notation.

9) Define DAX: Demonstrate any 10 DAX formulas with a context.

i) Count: Counts the no. of non-blank rows in a column. Typically used to count the no. of rows with valid entries in a dataset.

ii) Distinct count: Counts the distinct (unique) values in a column. Useful for identifying the no. of unique entries, such as distinct customer products or transactions.

iii) Average: Calculates the average of a column used to find the mean value, such as average sales or average quantities.

iv) 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.

v) 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. Table Page No. 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.

9) Sum: Sums the values in a column used to add up numerical values such as sales amount, quantities or profits.

10) 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.

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

Date: / /

Page No.:

Calculated Measure

- * Calculated on the fly based on context
- * Based on filter and slicer context
- * It not stored calculated dynamically
- * used in aggregations and summaries
- * Performance impact recalculation each time it's used
- * filter use cannot be used directly in filters
- * No impact on model structure
- * It is dynamic and changes with context.

Calculated Column

- * Row wise calculation stored in the model
- * calculated for each row, independent of filters
- * stored in the data model
- * used for row level calculations or new dimension
- * Calculated once stored at load time.
- * can be used directly in filters & slicers.
- * Adds new columns to the model
- * static once calculated.
- * Operates on individual rows
- * often use IF() ~~condition~~ ^{condition} etc.

* Operates on aggregates or summaries
* use func like Sum(), Average()