

## Principles of Data visualization

- 1) what is web scraping. Describe the steps to perform web scraping.

Web Scraping is the process of automatically extracting data from websites. It allows users to gather information from web pages and store it in a structured format for analysis such as in a database or a spreadsheet. Web Scraping is widely used in applications like data analysis, price comparison, lead generation, and research. The data is often scraped using automated tools and scripts, typically written in languages like Python, which help navigate and extract the required information from HTML pages.

### Steps to Perform Web Scraping:

- 1) Identify the target website: Choose a website from which data needs to be extracted. Analyze the structure and format of the website, ensuring that the data is publicly available and that scraping the site complies with its terms of service.
- 2) Inspect the web page: Use the browser's developer tools (Inspect) to examine the HTML structure

of the web page. Identify the specific tags (such as `<div>`, `<table>`, `<span>`, etc.) that contains the data you want to scrape.

for example in a e-commerce site, product name and prices may be stored in specific HTML elements.

### 3) Choose the Appropriate Tools or Libraries:

Depending on the programming languages, choose appropriate libraries or frameworks for web scraping. In python popular libraries include:

- \* [BeautifulSoup](#): for parsing HTML and extracting data.
- \* [selenium](#): for simulating browser interaction & scraping dynamic content.
- \* [Scrapy](#): A comprehensive framework for large scale scraping projects.

### 4) Send HTTP Requests:

use tools like requests to send HTTP GET Requests to the server and retrieve the HTML content of the web page. the response from the server contains the raw HTML data.

import requests

```
response = requests.get('https://example.com')
```

html\_content = response.content

### 5) Parse the HTML Content:

once the HTML content is fetched, parse it using a library like BeautifulSoup. This allows for easy navigation of the document.

ture & extraction of data from specific tags & attributes from bs4 Import beautifulsoup.

soup = BeautifulSoup(html-content, 'HTML parser')

titles = soup.find-all('h1')

- 6) Extract required data: use the HTML parser to locate & extract the relevant data, this might involve searching for specific tags, classes or IDs that contain the info.

prices = soup.find-all('span', class\_ = 'price')

- 7) Handle pagination: many website split data across multiple pages. If the target website uses pagination identify the URL patterns or navigation buttons & write logic to scrape all relevant pages iteratively.

- 8) Store the data: After extraction, store the data in a structured format, such as csv file, json file or database.

import CSV

with open('data.csv', 'w') as file:

writer = CSV.writer(file)

writer.writerow(['Title', 'Price'])

for title, price in zip(titles, prices):

writer.writerow([title.text, price.text])

- 9) Handle JavaScript content: some websites load content dynamically using Javascript. Tools like Selenium can be used to simulate a browser environment and extract data after Javascript execution

② Respect Ethical & Legal Boundaries: Always check the website's robots.txt file to understand the permissions for crawling or scraping. Ensure compliance with legal regulations, including copyright laws & terms of service, and avoid overloading servers with excessive requests (use delays between requests).

③ Challenges in Performing Web Scraping:  
Web scraping, while an effective method of extracting data from websites, comes with several challenges and limitations. These challenges can arise from legal, technical or ethical considerations, making it critical for developers to address them when performing scraping tasks.

#### 1) Legal and Ethical Issues:

- \* Terms Of Services (TOS)
- \* Copyright and Data ownership
- \* Robots.txt & crawling guidelines

#### 2) Dynamic Content and JavaScript Rendering:

- \* Javascript loaded content
- \* Handling dynamic pages

#### 3) Website Structure Changes

- \* frequent HTML Structure changes
- \* Inconsistent HTML Tags.

## 4) CAPTCHA &amp; Bot Detection:

- \* CAPTCHA systems \* Antiscraping measures
- \* IP Blocking & Rate limiting.

## 5) Data Quality Issues:

- \* Incomplete or Noisy data
- \* Encoding and formatting issues.

## 6) Ethical Considerations &amp; Server Overload:

- \* Server load & Bandwidth usage
- \* Respect for website bandwidth.

## 7) Data Reliability and Updates:

- \* Real-time data changes \* Consistency over time.

## 8) Data Privacy &amp; Regulations:

- \* Personal Data protection \* Regulatory compliance.

## B) Demonstrate a Python web scraping script using BeautifulSoup library for navigating any HTML.

A) BeautifulSoup is a Python library that allows easy navigation, searching and modification of HTML and XML documents. It is particularly useful for web scraping as it helps parse HTML files & extract the desired data.

```
from bs4 import BeautifulSoup
```

```
with open('example.html', 'r', encoding='utf-8') as file:
```

```
    html_content = file.read()
```

```
soup = BeautifulSoup(html_content, 'html.parser')
```

```
page_title = soup.title.text
```

```
print("page-title:", page-title)
books = soup.findAll('div', class_="book")
for book in books:
```

Page No.:

```
    title = book.find('h2').text
    print("Book title:", title)
    author = book.find('p').text
    print("author:", author)
    price = book.findAll('p')[1].text
    print('price:', price)
    print('--')
```

Output: Page title : Book store.

Book title = Python for Beginners

Author = John Doe

Price : \$30 .

v)

Illustrate the difference between loc & iloc function in python program with a code snippet;

A In python, loc & iloc are functions provided by the pandas library to access data from a Datapframe.

\* loc is used for label-based indexing. It allows you to select rows & columns by labels.

(row/column names or indices)

\* iloc is used for position-based indexing. It selects rows and columns by integer positions (row/column index numbers).

import pandas as pd

```
data = {'Name': ['Akash', 'Nithik', 'Karthik']}
```

'Age' : [22, 23, 26],  
 'Score' : [85 90 86] ]

```
df = pd.DataFrame(data, index=['a','b','c'])

print("Original DataFrame : ")
print("Using loc to select row b & columns 'a' ")
print(df.loc['b', ['Name', 'Age']])
print("Using iloc to select row (1) & first two columns")
print(df.iloc[1, [0, 1]])
print("Using loc to select multiple rows : ")
print(df.loc['a':'b', ['Name', 'Score']])
print("Using iloc to select multiple rows : ")
print(df.iloc[1:2, [0, 2]])
```

Output : Original DataFrame			using loc
	Name	Age	Score
a	Akash	22	85
b	Nithik	23	90
c	Karthik	26	86

  

using loc to select multiple rows			using iloc
	Name	Age	Score
a	Akash	22	85
c	Karthik	26	86

  

using iloc to select multiple rows			using iloc to select multiple rows
	Name	Age	Score
1	Nithik	23	90
2	Karthik	26	86

loc uses label to select data.

iloc uses integer position to select data.

Write a python program to delete 3rd column, delete  
3rd element & insert new element in 3rd column  
insert array in np.

array = np.array([1, 2, 3, 4, 5, 6])

(7, 8, 9, 10)

(7, 8, 9, 10)

print("Original array (array)")

print(array)

array = np.delete(array, 2, axis=1)

print("array without 3rd col.")

array = np.insert(array, 2, 4567)

array = np.insert(array, 2, newarray, axis=1)

print("array - new")

Original array. Array after deleting 3rd column

1 2 3

1 2 1

4 5 6

4 5 1

7 8 9

7 8 1

array after inserting in 3rd column.

1 2 10

4 5 11

7 8 12