

## FIFTH PROGRAM

web scraping is also known as web crawling. It is the process of extracting data from websites. web crawlers (spiders or bots) automatically navigate the web by following links from one page to another.

The main goal is to discover new pages and revisit existing ones to see if they have been updated.

Crawlers download the content of web pages for further processing.

Indexing is the process of processing and storing the data collected by web crawlers into a structured format. it is easy and fast to retrieve the data. ex: sports, movies, politics.

Sitemap is a file written in XML. It lists all the URLs of a website. Site owners submit Sitemap data. If you want to design a website. your website is able to show in the search engine you should give crawling and indexing.

Send HTTP Requests:- Using program language like python requests library to send request to web servers and retrieve the HTML content of web pages.

Parse HTML:- Use parsing libraries like BeautifulSoup to extract specific data from HTML (Heading, paragraph, links, tables, images).

Extract Data:- with the parsed HTML. you can locate and extract the relevant data using HTML tags.



\_/\_/\_

Cleaning and Structuring:- After extracting you need to clean & structure it.

Finally Store Data:- Save the extracted data in a structured format like CSV, JSON or database for further analysis.

Robots.txt file:- It tells search engine which pages to access and which are not access.

User-agent: \* All  
Command disallow → folder/pages.

User-agent: |directory|

The 'robots.txt' file must be placed in the root directory of the website.

robots.txt file consists of one or more blocks of directives.

each block starts with user-agent line to specify which crawlers the directives apply followed by one or more Disallow (or) Allow line

ex:-  
User-agent: \*  
Disallow: /private/  
Disallow: /tmp/  
Allow: /public/  
User-agent: Googlebot  
Disallow: /no-google.

# Block all web crawlers from accessing the private directory.

User-agent: \*  
Disallow: /private/



# Allow all crawlers to access the public directory  
Allow: /public/

# prevent the Google web crawler from accessing the /no-google directory.

User-agent: Googlebot

Disallow: /no-google/

# Provide the location of sitemap

Sitemap: https://www.sample.com/sitemap.xml

import requests

response = requests.get("http://en.wikipedia.org/robots.txt")

text = response.text

print("robots.txt for http://www.wikipedia.org/")

print(text)

(08)

import requests

# make a get request to the robots.txt file.

response = request.get("\_\_\_\_\_")

# Check if the request was successful.

if response.status\_code == 200:

# Store the content of the robots.txt file

text = response.text

# print the contents

print(text)

else:

print("Failed to retrieve the robots.txt file")