



Web Scraper Project Algorithmic Analysis **Pech Aque, Damaris Esther**

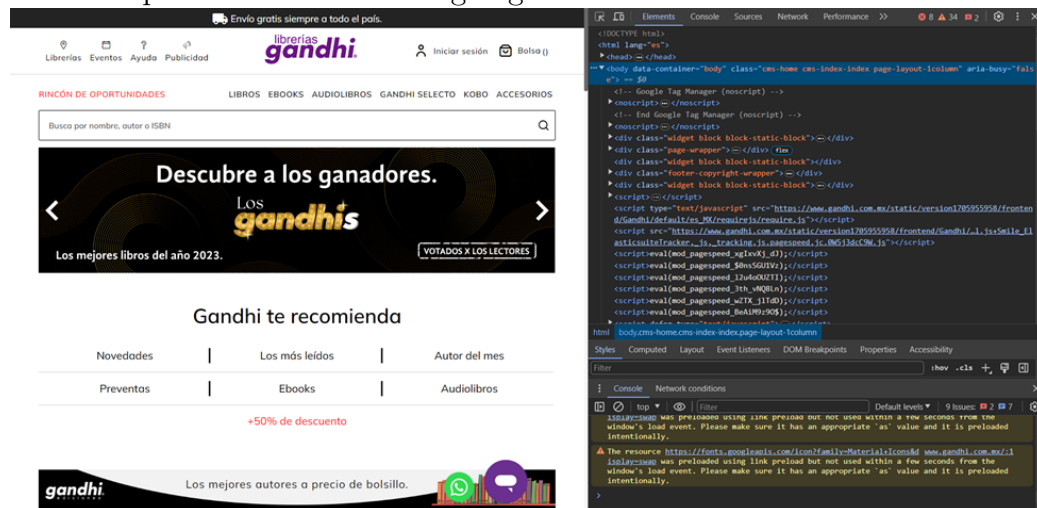
Lopez Tzec, Isaias de Jesus
Herrera Balam, Daniel Fernando
Paredes Dzib, Valeria de los Angeles
Velasco Martin, Jonathan Abisai

4th Year A - Data Engineering
Polytechnic University of Yucatán
Professor: Adrian Roberto Carmona Rodriguez
Date: January 22, 2024

Web Scraper

Preliminary Analysis of the Website

We start analyzing the website looking for what kind of data we are interested in, how the URL works for the research, and how the HTML is organized. Knowing all this information is going to make the development of the script easier to scrape all the data we are going to need.



Planning the Scraper

After analyzing the website, when you write a keyword for the research, a lot of information appears for books like name, author, price, kind of object (type of cover of the book).



Legal and Compliance Review

The web scraping project for the bookstore page underwent a comprehensive Legal and Compliance Review, ensuring adherence to ethical and legal standards.

Selection of Tools and Languages

We decided to use Python for the scraper because it is the base for our career and also one of the easiest languages for tasks like this.

Development of the Scraper

As mentioned, we use Python and a few libraries for our scraper. BS4 lets us work with the HTML of the page and manipulate it for what we want.

```

Ghandipy X
C:\Users\> this > Desktop > Carmona cuatri 4 > Ghandipy > ...

1 | Import of required libraries
2 | import pandas as pd
3 | from urllib.request import urlopen
4 | from bs4 import BeautifulSoup
5 | import requests
6 | import mysql.connector
7 |
8 | # Definition of headers to simulate a browser request
9 | headers = {
10 |     "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64; rv:94.0) Gecko/20100101 Firefox/94.0",
11 |     "Accept-Language": "en-PK, en;q=0.9"
12 | }
13 |
14 | # Definition of the search string and the list for storing books
15 | search_query = 'Harry'
16 | books=[]
17 |
18 | # Main loop to iterate through 10 pages of results
19 | for page in range(1,11):
20 |     # Construction of the search URL with page number and search string
21 |     search_url = f"https://www.gandhi.com.mx/catalogsearch/result/index/?p={page}&q={search_query}"
22 |     # Performing the GET request to the search URL with the headers defined
23 |     response = requests.get(search_url, headers=headers)

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

PS C:\Users\titlis > C:\Users\titlis\AppData\Local\Programs\Python\Python38\python.exe "C:\Users\titlis\Desktop\Carmona cuatri 4\Ghandi.py"

```

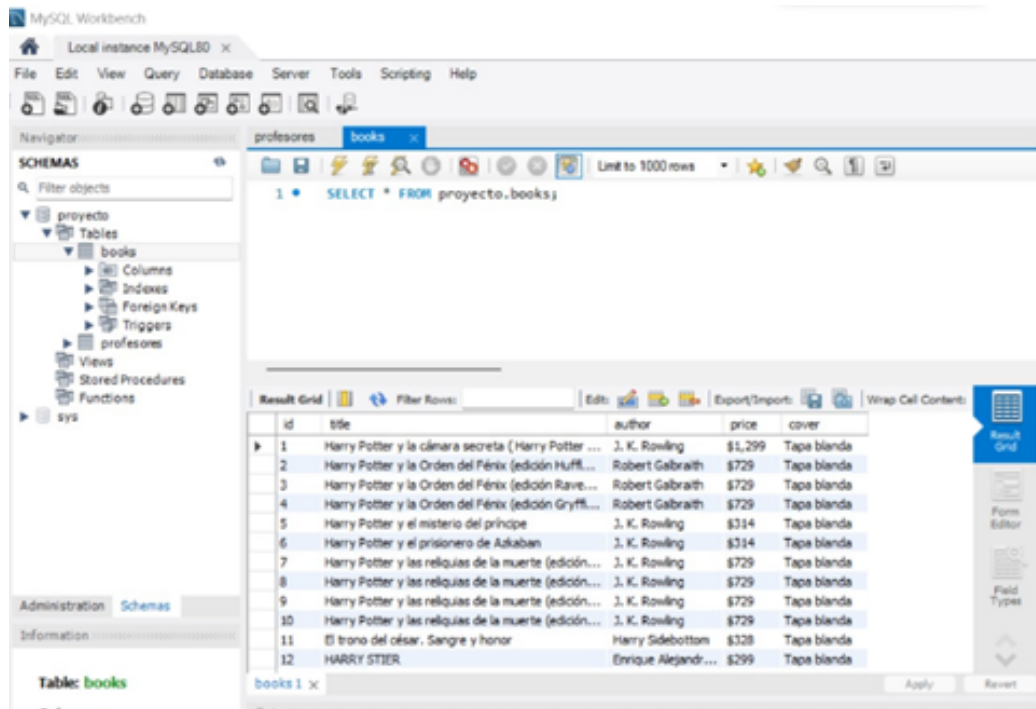
	Title	Author	Price	Cover
0	Harry Potter y la cámara secreta (Harry Potte...	J. K. Rowling	\$1,299	Tapa blanda
1	Harry Potter y la Orden del Fénix (edición ilu...	Robert Galbraith	\$729	Tapa blanda
2	Harry Potter y la Orden del Fénix (edición Row...	Robert Galbraith	\$729	Tapa blanda
3	Harry Potter y la Orden del Fénix (edición Gry...	Robert Galbraith	\$729	Tapa blanda
4	Harry Potter y el misterio del príncipe	J. K. Rowling	\$314	Tapa blanda
5	Harry Potter y el prisionero de Azkaban	J. K. Rowling	\$314	Tapa blanda
6	Harry Potter y las reliquias de la muerte (edi...	J. K. Rowling	\$729	Tapa blanda
7	Harry Potter y las reliquias de la muerte (edi...	J. K. Rowling	\$729	Tapa blanda
8	Harry Potter y las reliquias de la muerte (edi...	J. K. Rowling	\$729	Tapa blanda
9	Harry Potter y las reliquias de la muerte (edi...	J. K. Rowling	\$729	Tapa blanda

Efficiency and Performance Management

To enhance the efficiency and performance of the web scraping process, several measures were implemented.

Data Storage

For effective data storage, a structured approach was taken by utilizing a MySQL database.



Testing and Debugging

Thorough testing procedures were incorporated to ensure the correct collection of data during the scraping process.

Documentation and Maintenance

Comprehensive documentation practices were followed to provide clarity on both the code and the scraping process.

Implementation and Monitoring

For efficient deployment and monitoring, the web scraper is planned to be encapsulated within a Docker container.

