

Task Report



Website Information Extraction

Submitted

By

Anuroop Arya

2024

Acknowledgment

I want to express my heartfelt appreciation to AllHeartWeb for providing me with the chance to work on this project. Their project has really aided in my professional development in the field of data science and improved my grasp of web scraping and data extraction. I appreciate AllHeartWeb's unwavering support and giving me a stage on which to hone my technological abilities.

Project Overview

The objective of this project was to extract and analyze various attributes from a list of 100 websites, including social media links, tech stack information, meta tags (title and description), payment gateways, website language, and website category. The extracted data was stored in a MySQL database for further analysis.

Approach

Programming Language: Python

Tools and Libraries Used:

- `requests` for fetching HTML content
- `BeautifulSoup` for parsing HTML
- `langdetect` for language detection
- Regular expressions for extracting specific patterns

Data Extraction:

- HTML Content: Fetched from each URL using the `requests` library.
- Parsing: HTML content parsed using `BeautifulSoup` to extract relevant information.

Information Extraction:

- Social Media Links: Identified and extracted using regular expressions.
- Tech Stack: Keywords related to MVC frameworks, CMS, and JavaScript identified within the HTML content.
- Meta Tags: Title and description extracted directly from the HTML.
- Payment Gateways: Links to popular payment gateways identified using regular expressions.
- Website Language: Primary language identified using the `langdetect` library.
- Website Category: Classified using a predefined set of keywords.

Data Storage:

- Extracted data stored in a MySQL database named `website_info`.

Error Handling:

- Robust error handling implemented to manage exceptions during data extraction and insertion processes.

Challenges Faced and Solutions Implemented

1. Handling Dynamic Content:

Some websites loaded content dynamically using JavaScript, necessitating additional handling to fetch complete content. Addressed by ensuring `requests` library fetched complete rendered HTML content.

2. Website Variability:

Different websites had different structures and content, requiring adjustments in extraction patterns and rules. Flexible regular expressions and keyword searches were used to adapt to various website layouts.

3. Language Detection Issues:

Some websites contained multiple languages or non-standard language text, occasionally causing errors in language detection. Mitigated by implementing the `langdetect` library.

4. Database Operations:

Ensured data integrity and accurate insertion into MySQL database (`website_info`), particularly handling special characters and formatting issues. Used parameterized queries and proper encoding.

Learnings and Conclusion

This project, which involved the scraping of 100 websites, was my first foray into web scraping and data extraction. It offered insightful information about the intricacies of web data and useful uses for data science methods. Through this project, I improved my database management, web scraping, data manipulation, and Python programming skills.

In the future, I want to improve this project by researching more data properties for extraction, improving error handling procedures, and streamlining the database structure for more efficiency. My career as a data science intern has a strong foundation thanks to this project, which has given me the hands-on experience and abilities I need to take on increasingly difficult data challenges in the future.