

Project 1 Report

INTRODUCTION:

Objective:

The primary aim of this project is to leverage web scraping techniques to extract valuable data from the Yellow Pages website. This dataset will serve as a starting point for future research projects. It includes a complete list of restaurants in a specific place, including key insights such as restaurant names, contact information, precise addresses, customer ratings, and a variety of other relevant information.

Motivation:

This initiative is a useful resource for both individuals and businesses. It simplifies the process of selecting the appropriate restaurant for diners in Boston, MA by giving a full list complete with contact information, addresses, ratings, and amenities. Businesses can use this data to improve their services, personalize their products, and strategically position themselves for success in high-potential places at the same time. This dual-purpose effort attempts to improve the dining experience for consumers while also offering useful insights for food industry businesses.

Out of Scope:

While this project succeeds in collecting dynamic data from the Yellow Pages website, it's crucial to highlight that the current scope is mostly focused on the web scraping process. As a result, advanced data analytics, machine learning methods, and the development of comprehensive front-end interfaces are not within the scope of this project. The dataset developed in the project, on the other hand, provides as a good foundation for future projects that may include more intricate data analysis techniques or sophisticated applications.

METHODOLOGY:

Environment setup and tools used:

The code is implemented in Python, written on Jupyter Notebook and relies on several libraries:

- **requests:** For making HTTP requests to the Yellow Pages website.
- **BeautifulSoup:** For parsing the HTML content and extracting relevant information.
- **Selenium:** Used to handle dynamic content loading.
- **Pandas:** To covert output to dataframe

SOLUTION:

A sophisticated yet methodical procedure was developed as the answer to this web scraping project in order to guarantee the precise extraction of important restaurant data. The following steps are included in it:

Dynamic Content Handling with Selenium: I used Selenium, a potent automation tool, to load dynamically generated web pages. By doing so, I was able to interact with the page's elements and access data that would not have been instantly available in the HTML response.

Optimizing for Dynamic Content Loading: Built a deliberate wait mechanism, giving dynamic material enough time to fully load, to account for changes in loading times. This makes sure that all relevant data is available before beginning the extraction process.

Using BeautifulSoup to Extraction Precise Data: The parsing of the HTML material was crucial and relied heavily on the BeautifulSoup library. To browse the complex HTML structure and retrieve data pieces, such as restaurant names, phone numbers, addresses, websites, ratings, and amenities, it was used to its fullest potential.

Organized Data Structures: To display the retrieved information systematically as a dataframe, used Pandas.

Use of Extracted Data:

The name, telephone, address and website of the restaurant will give prospective consumers the basic information they need to quickly find and contact each establishment. This crucial information lays the foundation for a seamless and effective dining experience and empowers customers to choose meals sensibly. Additionally, this information provides firms with an essential point of contact for client connection, facilitating bookings, questions, and feedback.

Ratings are a useful metric for evaluating client preferences and satisfaction levels in a restaurant recommender system. These measurable ratings provide explicit feedback on the dining experience, allowing the system to create personalised recommendations based on customer preferences and expectations.

Amenities, which include additional services and conveniences, play an important part in improving the dining experience. Knowing the facilities provided by each restaurant is critical for a recommender system. It enables individualized suggestions, ensuring that consumers are paired with places that meet their specific likes and needs.

Handling HTML and CSS:

The HTML-parsing powerhouse BeautifulSoup package was used to effortlessly navigate and retrieve specific data from the HTML source. We made sure that the crucial restaurant information would be accurately retrieved by identifying classes and tags.

Expansion of Dynamic Content:

With the help of Selenium, this project is capable of handling dynamic content loading with ease. The code can be further extended to effectively handle many pages of results in the case of additional pagination needs, providing a complete dataset.

Rate Limiting and Error Handling:

A thorough error handling mechanism was put in place to ensure that the web scraping process was resilient. To ensure that each request is successful, this includes thorough examinations of the HTTP response status codes. In the event of an error, a clear and informative message is generated, aiding in rapid issue resolution.

Data Storage and Export:

The data-frame output can be stored and exported as CSV file and can be used for future analysis.

CHALLENGES & OUTLOOK:

Challenges:

- 1) When I tried to use Selenium, was getting issue in executing it, then tried installing suitable chrome drivers which worked fine.
- 2) The dataframe was not displaying in one page, last 2 columns were going to next line even if there was space in the previous page. So, used maximum columns and maximum colwidth functions.

Future Improvements:

- 1) Data cleaning can be implemented to remove empty/None/NaN values and in amenities columns, all the available amenities do not have spaces in between, so suitable changes can be implemented for better display.
- 2) Can work in depth for pagination for multiple pages result.