# Week 3

## What is Web Scraping?

Web scraping is an automatic method used to efficiently obtain a large amount of data from websites. Its primary function is converting unstructured HTML or other types of data into a structured data format that is easier to analyze and store.

For instance, when large e-commerce sites like Amazon or Flipkart contain massive amounts of data, web scraping is the technique employed to systematically extract this information.

## The Web Scraping Process

The flow of data through a web scraping operation involves three main stages:

1. Websites: This is the source, providing the original unstructured data.

2. Scraping Platform: This tool or framework is responsible for processing the data request, navigating the site, and performing the extraction.

3. Structured Data: This is the final output, organized into a usable format, such as a table, a database, or a CSV file.

---

## Applications and Strategic Benefits

The structured data obtained through web scraping is highly versatile. It can be used for fundamental data analysis, developing predictive models, and even for populating content on other websites.

## Key Applications

Web scraping is a cornerstone technology in several data-driven fields:

- Price Monitoring: Tracking competitor pricing in real-time.

- Market Research: Gathering product and consumer information to identify trends.

- News Monitoring: Collecting and aggregating articles and updates from various news sources.

- Sentimental Analysis: Extracting public opinions and sentiment from forums, social media, or reviews.

- Email Marketing: Collecting contact information for outreach.

## Why Companies Use Scraping?

Companies leverage web scraping for significant strategic and analytical advantages:

- Optimal Pricing: Firms can scrape both their own product data and that of competing products. This comparison is critical for understanding market impact and setting optimal pricing strategies.

- Consumer Trend Analysis: Scraping helps in analyzing current consumer trends, which is essential for determining the company's future direction and strategy.

- Technical Data Extraction: It enables the extraction of specific technical details, such as software versions, end-of-life (EOL) dates, release years, product names, IP addresses, and hostnames.

---

**Web Scraping Tools: Focus on Selenium**

Web scraping can be executed using various programming tools. In the Python ecosystem, popular choices include libraries like BeautifulSoup and Selenium, as well as full-fledged frameworks like Scrapy. This document specifically details Selenium.

**What is Selenium?**

Selenium is a robust tool designed for controlling web browsers programmatically, effectively enabling browser automation.

- It is highly compatible, functioning across all major web browsers and supporting all major Operating Systems (OS).

- Scripts can be developed in multiple programming languages, notably Python and Java.

- A fundamental component of the tool is the WebDriver, which acts as an API (Application Programming Interface), providing the necessary classes, functions, and methods to interact with the browser.

**Selenium WebDriver Features**

The Selenium WebDriver is utilized for automation due to its comprehensive features:

- Multi-Browser Compatibility

- Multiple Language Support

- Speed & Performance

- Community Support

- Open Source & Portable

- Work on Different OS

- Add-ons & Reusability

- Simple Commands

- Reduced test execution time

- No server installation

---

**Selenium Architecture and Implementation**

Selenium Architecture

The architecture is designed to handle the interaction between the user's code and the actual web browser:

1. Selenium Client Libraries: These are the pre-developed libraries for different languages (like Python, Java) that allow the user to write their automation code.

2. JSON Wire Protocol: This protocol serves as the intermediate translator. It takes the client's request and converts it into a format understandable by the Browser Driver (the server) and also translates the server's response back to the client.

3. Browser Drivers: These specialized components (e.g., ChromeDriver, FirefoxDriver) establish a secure connection with their respective browsers, executing commands without exposing the browser's underlying logic.

4. Real Browser: This is the actual web browser instance (e.g., Chrome, Firefox, Microsoft Edge) that receives the commands and renders the webpage.

**Implementation Steps for Scraping**

To begin a web scraping project using Selenium, follow these general steps:

1. Installation: The Selenium library is installed using the Python package manager: pip install selenium.

2. Driver Requirement: Selenium requires a specific driver to interface with the chosen browser. The various available drivers include:

   o ChromeDriver

   o FirefoxDriver

   o InternetExplorerDriver

   o EdgeDriver

   o RemoteWebDriver

3. Driver Setup: An appropriate WebDriver (e.g., for Chrome) that is compatible with the local browser version must be installed and configured.

4. Import Libraries: Essential components such as webdriver, Service, By, and ChromeDriverManager are imported from the Selenium modules.

5. Driver Creation: The driver instance is created, often requiring the appropriate path to the browser driver executable.

6. Open Website: Navigation to the target URL is achieved using the command: driver.get("url").

7. Element Identification: The most critical step for extraction is locating elements on the webpage, typically done by specifying their XPATH (XML Path).

   o XPATH is a powerful language used for locating nodes within an XML (or HTML) document. It provides the precise path of the element on the webpage.

   o The basic syntax for locating elements by tag and attribute is: //tag_name[@attribute='value'].