

# Web Scraping

<sup>1</sup>Alok Singh, <sup>2</sup>Namrata Dhandra

<sup>1</sup>Scholar, <sup>2</sup>Professor

<sup>1</sup>Amity School of Engineering and Technology

<sup>1</sup>Amity University Uttar Pradesh, Lucknow, India

**Abstract :** The process of retrieving or extracting information from various websites is called “Web-Scraping”. As you know search engine is a mixed combination of computer hardware and software supplied by a company through which the website has been determined. It automatically collects all the information from the web through web crawlers that crawls the whole web periodically. Scraping the Search engines is the old trend since the internet got started. A great issue why search engines don’t want to do scraping with it is that they only wanted you to browse like a normal person.

**IndexTerms – Web Scraping, Selenium, Boilerpipe, Nutch, Celerity**

## I. INTRODUCTION

Data plays an important role in many field including marketing, scientific use, academic etc. Any one collects the data from various websites for their improvement. Copping of data from any website to any local machine is forbidden by almost every website authority for this reason a user manually copies the data from website to file storage. But this task is time consuming. For such reason web scraping is introduced. Web Scraping means accessing information from a website without stealing it. There are various reasons why search engines don’t give permissions to scrape. Many of the big companies like Google, the big dog, acts as a slow down websites but all of us know that what they want is to not access their data by any of us. Many people use proxies to hide the IP address so that they can do their work in very simple and easy way. IP address recognize people as a scraper when they takes help of proxies. Mostly people who scrapes the websites used it as a safety. A web crawler is basically a script of a program that browse the World Wide Web in an easiest manner.

## II. OVERVIEW OF WEB SCRAPING

It is a technique used for extracting data (Unstructured form) from the various websites and changing that data into simplified manner. It can also be identified as data extraction from web, web harvesting or screen scraping. Web scraping is a type of data mining. The basic aim of the this process is to extract information from a different websites and then transform the acquired data into structure like spreadsheets, database etc. That can be furthure manipulated. Extracting information from websites helps in decision making in business. As we all know, scraping technique is used to extract information from various web pages. Web pages are written in Hypertext Markup Language (HTML) or XHTML (based on XML). These web documents are represented using Document Object Model, (DOM tree) Here the aim of HTML is to provide the format of text to be displayed by browsers. From the operational point of view, a web scraping can be viewed as manual copy and paste. The difference is that it is done automatically, by a virtual agent. An agent is performing the same operation as a human when interacting with a web site. This agent can follow links, submit forms, browse web pages etc. Lets take an example If someone uses web scraping without considering policy for limiting requests, the server may find Denial-of-Service attack, due heavy request triggered within a short span of time[3]. In the next step parser goes through a user-specified paths within the document. These paths may be specified by CSS selectors. Generally such web scraping uses regular expressions to trim the information.

## III. WEB SCRAPING USES

Web Scrapers are used by Online Marketers to get data secretly from the various websites such as keywords, links, emails etc. Following are some areas:

- Change detection on website
- Product Price comparison
- Weather broadcasting and data monitoring
- Research analytics
- Analyze data in graphics
- Web Indexing & rank checking
- Advertisement analysis
- Market Analysis

#### IV. TECHNIQUES USED BY WEB SCRAPING

**Selenium** - Selenium is an internet browser device that has the capabilities to complete an extensive variety of assignments on autopilot. Figuring out how to use selenium will help you in seeing how exactly a website works. Selenium can assist you significantly something other than web scraping, such as testing sites and computerizing. So, using Selenium can make you a web scraping master.

**Boilerpipe** - While scraping clean content alongside related titles is the necessity, Boilerpipe is an incredible choice. BoilerPipe is a Java library made only to remove information from site pages. It can keenly evacuate pointless html labels found on the pages. The feature of Boilerpipe is that it can separate applicable substance in a matter of milliseconds and with negligible contribution from the client. The precision is stunningly high, which makes it one of the most straightforward apparatuses to use for scraping information. Getting acquainted with this apparatus can improve your web scraping aptitudes, immediately.

**Nutch** - Nutch is examined as the best quality level of web extracting advances. It is only an open source web crawler program that can slither and remove information from site pages at lightning speeds. Nutch can be utilized for creeping, fetching and putting away the information once customized for the particular prerequisite. To extract data from web, the website pages to cross and concentrate information from must be coded into Nutch physically. Static and dynamic pages can be recovered by presenting HTTP asks for on the remote web server utilizing attachment programming.

**Celerity** - Celerity is a great JRuby wrapper made around HtmlUnit – a Java program without using the head with help for JavaScript. It has a simple to utilize API that can be utilized to automatically explore through web applications. It is astonishingly quick since there is no tedious GUI rendering or pointless downloads. Being adaptable and non-meddlesome, it can keep running out of sight quietly after the underlying setup. Celerity is an awesome program computerization instrument you can use to rub the web effectively and quick.

#### V. SCRAPING A SITE

The primary thing you need to do while at the same time making a sitemap is deciding the start url. This is the url from where the scraping can takes place. You can in like manner show diverse to begin urls if the scratching will begin from various spots. For example in case you have to scrape diverse recorded records then you could make an alternate start url for each inquiry thing. A few sites looks at specific headers to be available and an exposed twist. You may not move beyond the captcha by and large (If there is any). There might be rate constrain applied (How numerous demand you can send in a given time span to a specific URL), so a scrapper content would effectively trigger that and site will quit reacting for that IP Address. Regardless of whether you get around this by utilizing different IP address, there is an opportunity to be delegated DDoS assault. Pages behind login may not be rejected utilizing basic contents. You will require a honest to goodness record and some kind of headless program and a testing structure to cooperate with the program utilizing code. One such framework is Google Chrome. To wrap things up, If a site doesn't enable the information to be rejected lawfully then you should forgo doing as such. With every one of these focuses here are a few focuses to help you in rejecting site. Attempt a basic twist or some other http library demand to some imperative pages of site which you think about most. On the off chance that that works then any content can parse and rub the site regardless of if its Python. Anyway Python's BeautifulSoup library is great. On the off chance that that doesn't work at that point attempt to look at the headers being sent by the real program when you visit the site typically and endeavor to emulate them utilizing twist. Check if there is any captcha and climate or not it can be skirted. Endeavor to make a major number of solicitations inside a limited capacity to focus time and check if there is any rate restrain mistake. In the event that yes you have to rest your content for the time indicated in rate confine headers and resume once more. Compose a basic content to simply download the HTML from most critical pages. On the off chance that that works then you can most likely rub the site. Numerous information investigation, huge information, and machine learning ventures require scraping sites to assemble the information on which you are working on. The python language is broadly utilized as a part of the data science network, and in this manner has a biological community of modules and apparatuses that you can use in your own venture.

#### VI. COMPONENTS OF WEBPAGE

When we visit a site page, our internet browser makes a demand to a web server. This ask for is known as a GET ask for, since we're getting records from the server. The server at that point sends back records that advise our program how to render the page for us. The documents fall into a couple of primary kinds:

- HTML — contain the fundamental substance of the page
- CSS — add styling to influence the page to look more pleasant.
- JS — Javascript records add intelligence to website pages.

After our program gets every one of the records, it renders the page and shows it to us. There's a great deal that occurs in the background to render a page pleasantly, yet we don't have to stress over the vast majority of it when we're web scraping. When we perform web scraping, we're occupied with the principle substance of the site page, so we take a gander at the HTML.

Hyper Text Markup Language (HTML) is a lang. that website pages are made in. HTML isn't a programming language, similar to Python — rather, it's a markup lang. that advises a program how to design content. HTML enables you to do comparative things to what you do in a word processor like Microsoft Word — influence content intense, to make para, etc. Since HTML isn't a programming dialect, it isn't so complicated as Python.

We should take a brisk visit through HTML so we know enough to scrape successfully. HTML comprises of components called labels. The most fundamental tag is the `<html>` tag. This label tells the internet browser that everything within it is HTML.

We haven't added any substance to our page yet, so in the event that we saw our HTML record in an internet browser, we wouldn't see anything:

We'll presently add our first substance to the page, as the `p` tag. The `p` tag characterizes a section, and any content inside the tag is appeared as a different para.

`a` and `p` are amazingly regular html labels. Here are a couple of others:

- `div` — demonstrates a division, or zone, of the page.
- `b` — bolds any content inside.
- `I` — emphasizes any content inside.
- `table` — makes a table.
- `shape` — makes an input form.

Class and id are the exceptional properties that give HTML components names, and make them simpler to interface with when we're scratching. One component can have different classes, and a class can be shared between components. Every component can just have one id, and an id must be utilized once on a page. Classes and ids are discretionary, and not all components will have them.

Beautiful Soup - Beautiful Soup is a Python bundle for parsing HTML and XML records (counting having deformed markup, i.e. non-shut labels, so named after label soup). It makes a parse tree for parsed pages that can be utilized to separate information from HTML, which is helpful for web scraping. Beautiful Soup gives a couple of basic strategies and Pythonic figures of speech for exploring, looking, and changing a parse tree: a toolbox for dismembering an archive and separating what you require. It doesn't take much code to compose an application. Beautiful Soup consequently changes over approaching archives to Unicode and active reports to UTF-8. You don't need to consider encodings, except if the archive doesn't determine an encoding and Beautiful Soup can't auto detect one.

## V. CONCLUSION

We should check a site's Terms and Conditions before you scrape it. As a rule, the information you extract ought not be utilized for business purposes. Do not ask data from the site too forcefully with your program (otherwise called spamming), as this may break the site. Ensure your program carries on in a sensible way (i.e. acts like a human). One requirement for one page for each second is great practice. The design of a site may change every once in a while, so try to return to the site and rework your code as required. Ideally inside a html tag, we put two different labels, the head tag, and the body tag. The principle substance of the website page goes into the body tag. The head tag contains information about the title of the page, and other data that for the most part isn't valuable in web scraping: Regardless, we haven't added any substance to our page (that goes inside the body tag), so we again won't see anything: In HTML, labels are settled, and can go inside different labels.

## REFERENCES

- [1] Ali, A. 2001. Macroeconomic variables as common pervasive risk factors and the empirical content of the Arbitrage Pricing Theory. *Journal of Empirical finance*, 5(3): 221–240.
- [2] Basu, S. 1997. The Investment Performance of Common Stocks in Relation to their Price to Earnings Ratio: A Test of the Efficient Markets Hypothesis. *Journal of Finance*, 33(3): 663-682.
- [3] Bhatti, U. and Hanif. M. 2010. Validity of Capital Assets Pricing Model. Evidence from KSE-Pakistan. *European Journal of Economics, Finance and Administrative Science*, 3 (20).
- [4] S.C.M. de S Sirisuriya, 2015, A Comparative Study on Web Scraping .Proceedings of 8th International Research Conference, KDU.
- [5] List of Web Harvester, Data Scraper, Web Scraping Software and Tools, n.d. WebData Scraping. URL <http://webdata-scraping.com/webscraping-software/>

- [6] Felipe Jordão Almeida Prado Mattosinho, Master Thesis, Mining Product Opinions and Reviews on the Web, TU Dresden.
- [7] Eloisa Vargiu, Mirko Urru1, 2013, Exploiting web scraping in a collaborative filtering- based approach to web advertising, Artificial Intelligence Research, 2013, Vol. 2, No. 1, <http://dx.doi.org/10.5430/air.v2n1p44>.
- [8] Schrenk, M. Webbots, spiders, and screen scrapers: a guide to developing Internet agents with PHP/CURL. No Starch Press, 2007.
- [9] OsmarCastrillo-Fernández,2015,WebScraping:Applications and Tools, European Public Sector Information Platform, Topic Report No. 2015/10.
- [10] Deepak Kumar Mahto, Lisha Singh, A Dive into Web Scraper World, 2016 InternationalConference on Computing for SustainableGlobal Development (INDIACom), 978-9-3805-4421-2/16/\$31.00 c , 2016 IEEE.
- [11] MiloslavBeio, Jakub Misek, Filip Zavoral, AgentMat:Framework for Data Scraping and Semantization, 9781-4244-2865-6/09/\$25.00 ©2009 IEEE.
- [12] Amir Ghazvinian, Sean Holbert, Nikil Viswanathan, Scrapy: Simple Web Scraping, Department of Biomedical Informatics, Department of Computer Science, Stanford University.

