



# Web Scrapping by Python feat. Beautiful Soup

Marcos R. Pesante Colón

ACM – CSE Chapter

August 8, 2020



# Web Scraping

Introduction & Applications

# Web Scraping

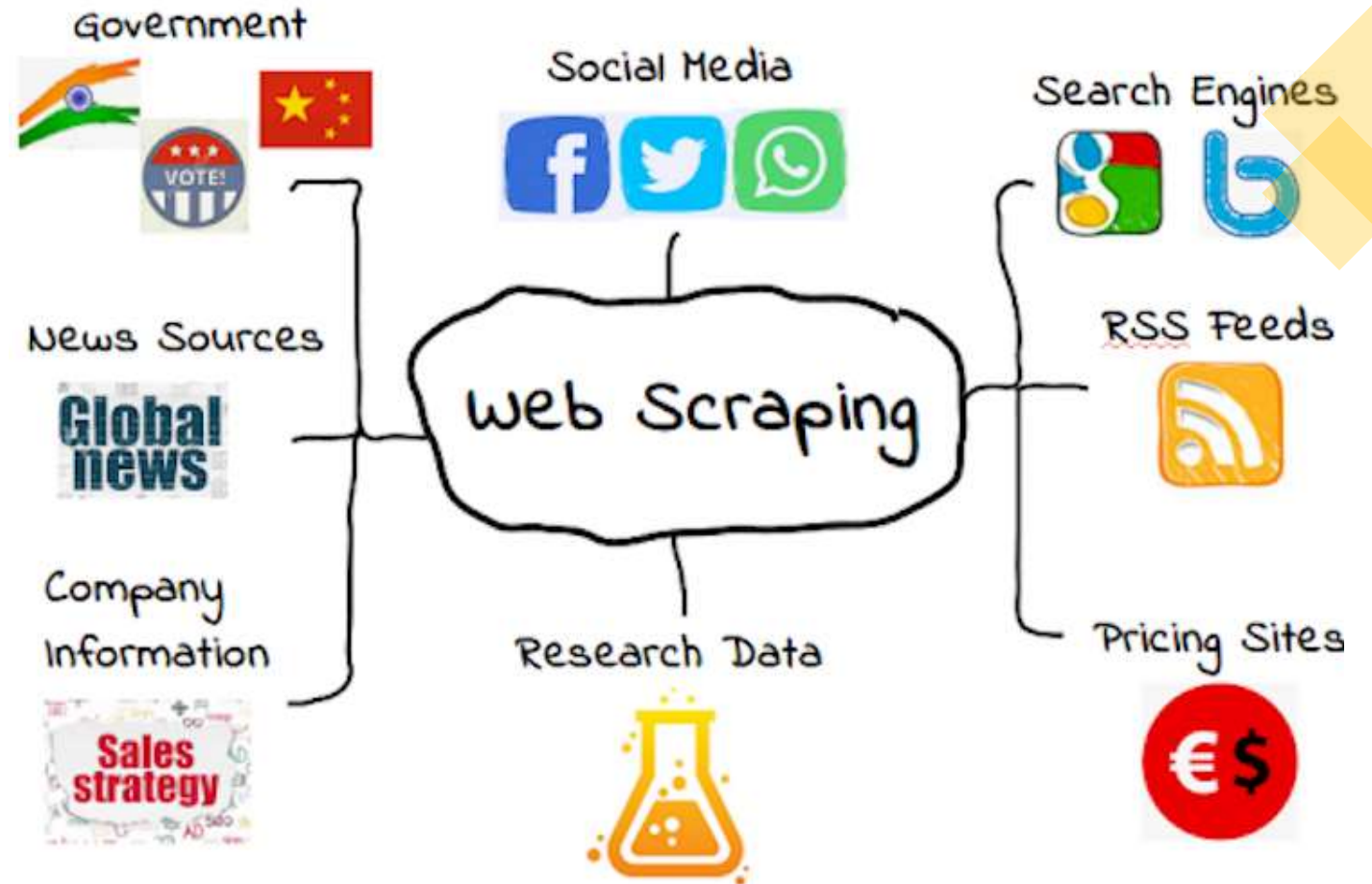
The extraction and copying of data from a website into a structured format using a computer program

<https://www.dictionary.com/browse/web-scraping>

\* Useful for when you want to automate data extraction from a website and they do not provide APIs or other services with this object

# Applications

- Google Web Crawling to find new websites
- Social Media Sentiment Analysis
- Price monitoring
- Obtain Data to feed Machine Learning Models



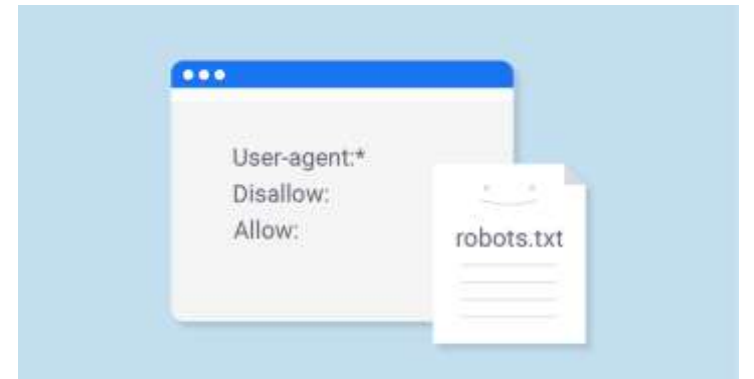
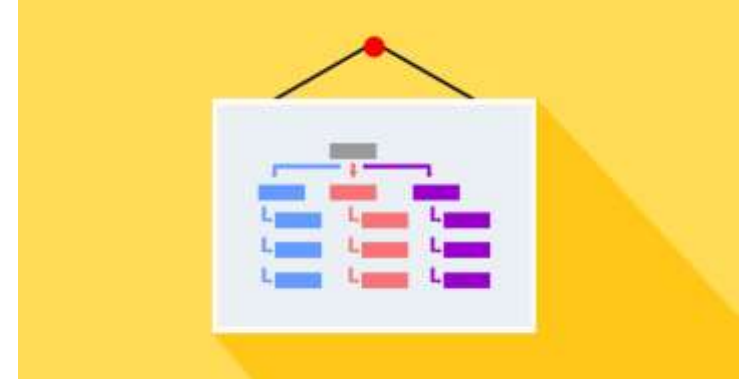
# Before Beginning

- Establishing goals
  1. What information do you want to extract?
  2. Who is going to use this data and how do they want it to be stored?
  3. How frequently will the data be “reset” ?



# Before Beginning

- robots.txt
  - List of restricted parts of the website
  - You cannot scrape these parts of the website
    - <https://moz.com/learn/seo/robotstxt>
- sitemap.xml
  - List of most, if not all, of the links in a website
  - Sometimes can even refer to even more sitemaps





# HTML

Introduction / Review

# HTML

- Hyper Text Markup Language
- Standard language to create and **structure** Web Pages
- Accomplishes this through the use of nested tags
  - `<h1></h1>`
  - `<p></p>`
  - `<img>`
  - `<a href="#"></a>`





# Ejemplo

```
<!DOCTYPE html>
<html lang="en">
  <head>
    <meta charset="utf-8" />
    <title>Page Title</title>
  </head>
  <body>
    <h1>This is a Heading</h1>
    <p>This is a paragraph.</p>
    <p>This is another paragraph.</p>
  </body>
</html>
```

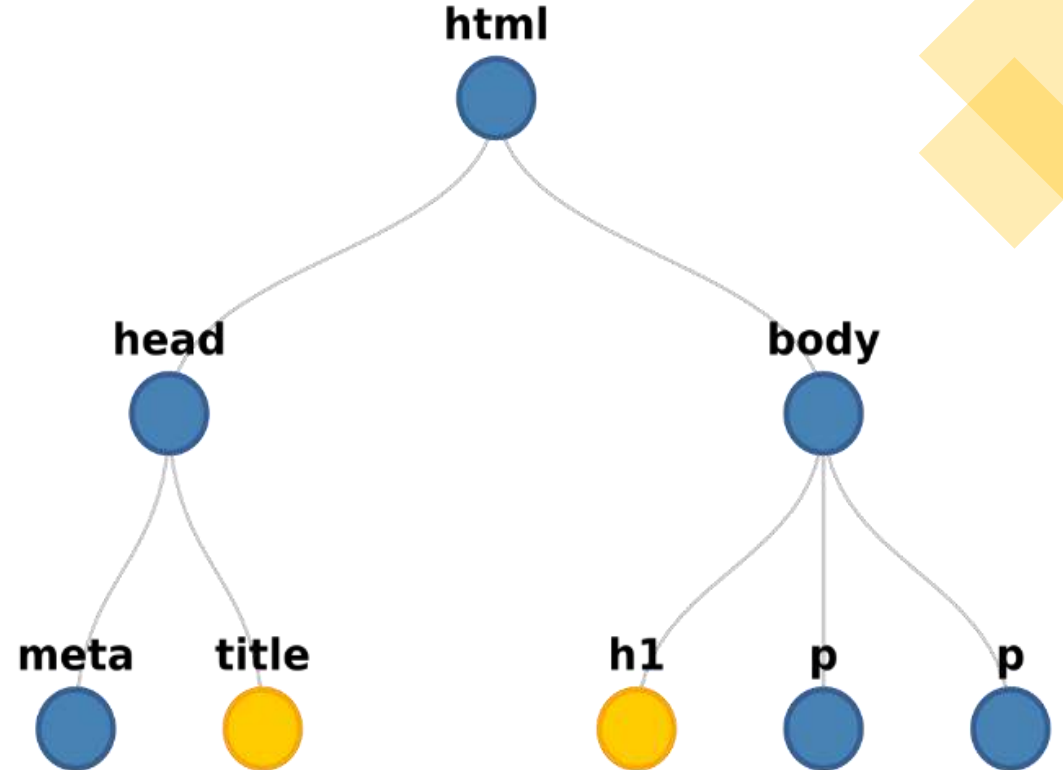
## This is a Heading

This is a paragraph.

This is another paragraph.

# Ejemplo

```
<!DOCTYPE html>
<html lang="en">
  <head>
    <meta charset="utf-8" />
    <title>Page Title</title>
  </head>
  <body>
    <h1>This is a Heading</h1>
    <p>This is a paragraph.</p>
    <p>This is another paragraph.</p>
  </body>
</html>
```



# Being picky about our HTML

CSS Selectors

# Where Python Comes In

Using Requests and  
Beautiful Soup



BeautifulSoup

