

Sistemas Web

Grado en Informática de Gestión
y Sistemas de Información

Dpto. de Ingeniería de Sistemas y Automática

Web Scrapping

Oskar Casquero (oskar.casquero@ehu.eus)
María Luz Álvarez (marialuz.alvarez@ehu.eus)



Esta obra está bajo una [licencia de Creative Commons Reconocimiento-CompartirIgual 4.0 Internacional](https://creativecommons.org/licenses/by-sa/4.0/).



Universidad
del País Vasco

Euskal Herriko
Unibertsitatea

BILBOKO
INGENIARITZA
ESKOLA
ESCUELA
DE INGENIERÍA
DE BILBAO



Estructura del tema

- Estructura de una página HTML
- Parsear una página HTML con BeautifulSoup
- Página descargada vs página renderizada
- *View Source Chart*: bookmarklet para ver la página renderizada
- *Selenium+Geckodriver*: automatizar navegador desde Python para renderizar HTML

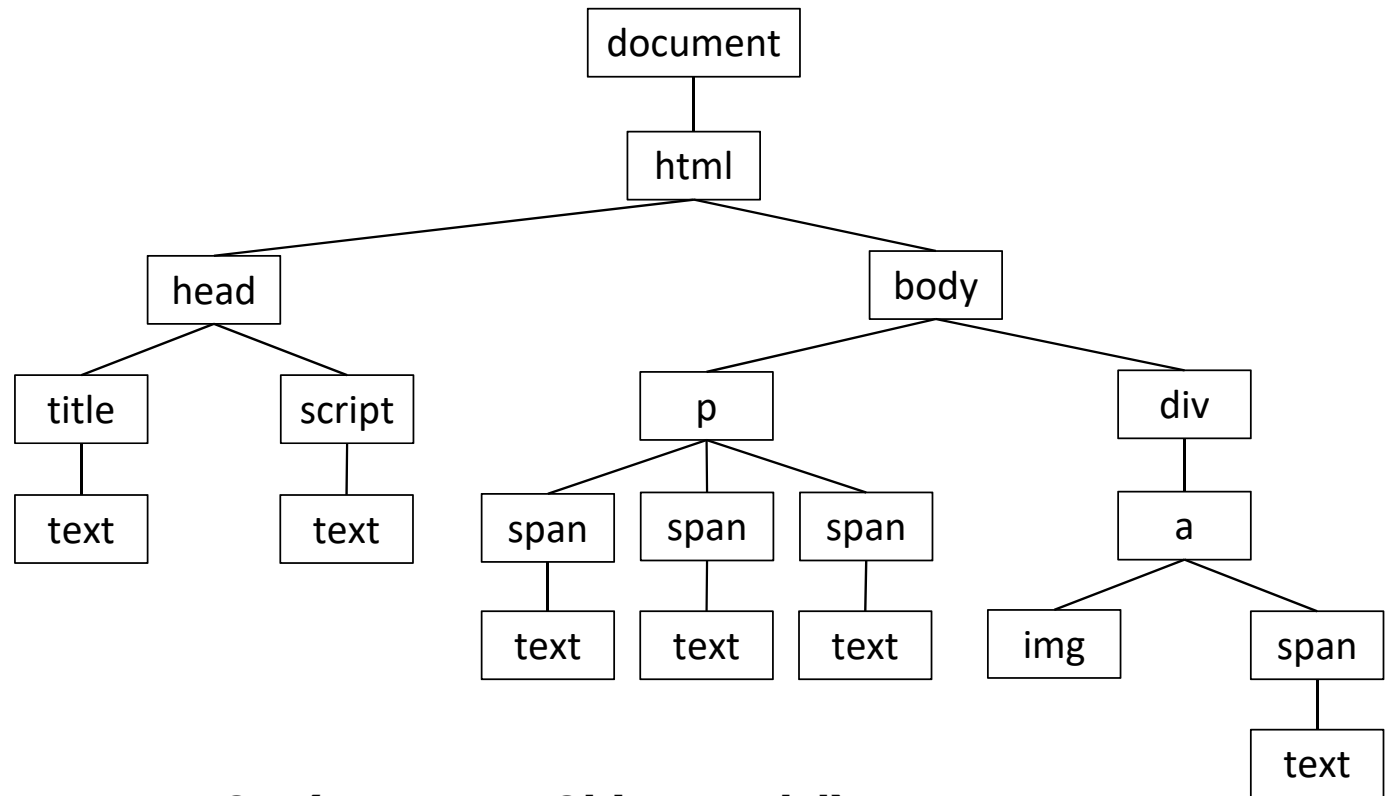
- **Caso de Estudio:** descargar imágenes de Google Images

- ☐ ¿La página HTML descargada coincide con la renderizada?
 - Sí → Obtener código HTML con *requests*
 - No → Obtener código HTML con *selenium+geckodriver*
- ☐ Analizar estructura de la página HTML para saber dónde están los enlaces a las imágenes.
- ☐ Parsear el HTML con BeautifulSoup
- ☐ Almacenar imágenes:
 - Realizamos peticiones HTTP para descargar las imágenes enlazadas
 - Decodificamos las imágenes “inline” de base64 a binario



1 Estructura de una página HTML

```
<html>
  <head>
    <title>título</title>
    <script>
      var data = new Date();
      document.write("Client Date: ");
      document.write(data);
    </script>
  </head>
  <body>
    <p>
      <span class="a">bla</span>
      <span class="a">bla</span>
      <span class="a">bla</span>
    </p>
    <div id="identifier">
      <a href="http://...">
        
        <span>bla</span>
      </a>
    </div>
  </body>
</html>
```



DOM (Document Object Model)

2 Parsear una página HTML

- Se utiliza la librería **Beautiful Soup** de Python para parsear HTML.
- Instalación de la librería: `pip install BeautifulSoup4`
- Documentación: <https://www.crummy.com/software/BeautifulSoup/bs4/doc/>
- **Ejercicio:** Programar un cliente que permita realizar consultas en el directorio de la UPV/EHU para mostrar un resumen con el nombre, apellidos y enlace de los resultados.

The image shows a web browser window on the left displaying the 'BILATU' directory page of the UPV/EHU. The page title is 'BILATU 2.4.2' and the URL is 'https://www.ehu.es/bilatu/buscar/'. The page content includes a search bar, a list of search results, and a sidebar with navigation links. The search results are titled 'RESULTADOS DE LA BÚSQUEDA' and list several names and their corresponding URLs.

On the right, a code editor window shows the Python script 'sending_form_directorio_beautiful_es.py'. The script uses the 'BeautifulSoup' library to parse the HTML of the search results page. The code includes comments and variable assignments for the search results, such as 'cuerpo = {\"abi_size\": sys.argv[1]}\" and 'respuesta = requests.request(metodo, uri, headers=cabeceras, data=cuerpo, allow_redirects=True)\".

3 Página descargada vs página renderizada

El navegador ejecuta el código javascript antes de renderizar la página.
El código javascript tiene capacidad para modificar la estructura y el contenido del HTML.

Página descargada

```
...  
<body>  
  Server Date: Fri Mar 04 08:43:33 CET 2016  
  <br/>  
  <script language="javascript">  
    var data = new Date();  
    document.write("Client Date: ");  
    document.write(data);  
  </script>  
</body>
```

**Abrir página
en Firefox**

Página renderizada

```
Server Date: Fri Mar 04 08:43:33 CET 2016  
Client Date: Fri Mar 04 08:43:34 GMT+0100
```

**Crear una página HTML
con este contenido**

4

View Source Chart: bookmarklet para ver la página renderizada

- La opción “Ver código fuente” del navegador permite ver el código de la página descargada.
- El bookmarklet “View Source Chart” permite ver el código de la página renderizada:
<http://viewsourcechart.com/getthebookmarklet.html> → **instalarlo en Firefox**

“View Source Chart”

“View Page Source”

```

1 <html>
2 <head>
3   <title>ERLAITZA</title>
4 </head>
5 <body>
6   <ul id="zerrenda">
7     <li>Server Date: Tue Feb 21 9:28:33 CET 2020</li>
8   </ul>
9   <br/> <!-- CRLF bezela -->
10  <script language="javascript">
11    var data = new Date();
12    var li_element = document.createElement("li");
13    var text_element = document.createTextNode("Client Date: " + data);
14    li_element.appendChild(text_element);
15    document.getElementById("zerrenda").appendChild(li_element);
16  </script>
17 </body>
18 </html><script>

```

```

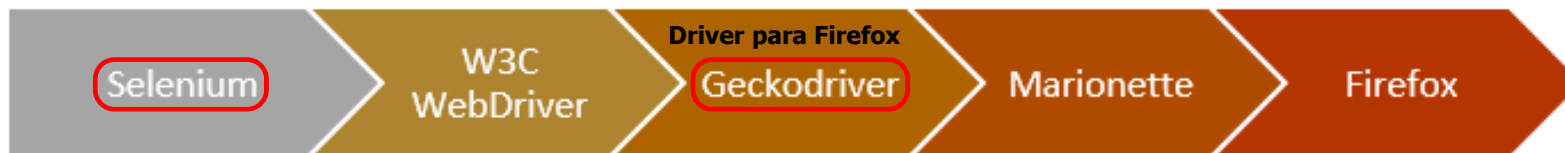
<html>
<head>
  <title>
    ERLAITZA
  </title>
</head>
<body>
  <ul id="zerrenda">
    <li>
      Server Date: Tue Feb 21 9:28:33 CET 2020
    </li>
    <li>
      Client Date: Fri Feb 18 2022 15:50:09 GMT+0100 (Central European Standard Time)
    </li>
  </ul>
  <br>
  <!-- CRLF bezela -->
  <script language="javascript">
    var data = new Date();
    var li_element = document.createElement("li");
    var text_element = document.createTextNode("Client Date: " + data);
    li_element.appendChild(text_element);
    document.getElementById("zerrenda").appendChild(li_element);
  </script>
</script>

```

Selenium+Geckodriver:

Automatizar navegador desde Python para renderizar HTML

- ¿Cómo podemos hacer uso del navegador desde un programa escrito en Python?
 - Mediante Selenium, un herramienta para la automatización de tests de aplicaciones web.
 - Geckodriver es un drive que ofrece una interfaz para controlar el navegador



```
from selenium import webdriver
from selenium.webdriver.support.ui import WebDriverWait
from selenium.webdriver.support import expected_conditions as EC
from selenium.webdriver.common.by import By

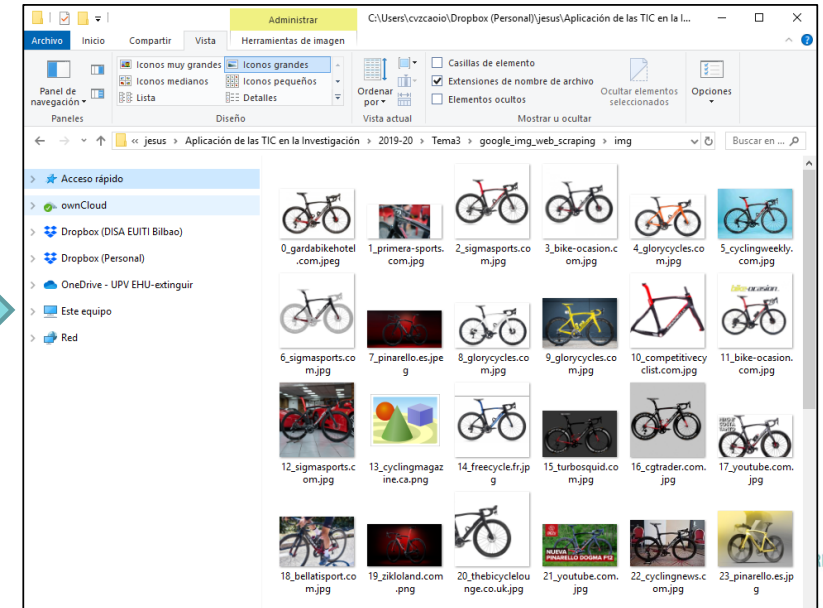
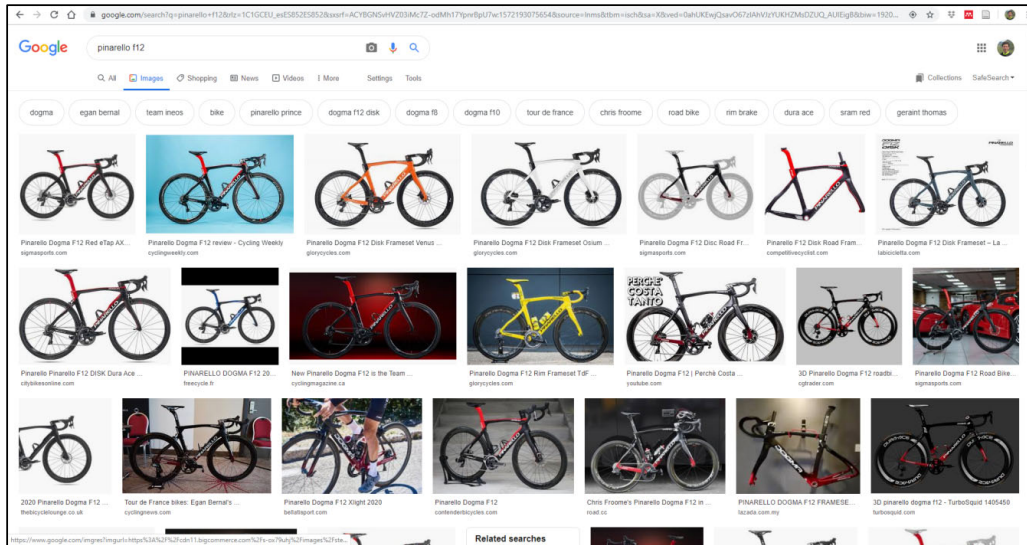
uri = "https://www.google.com/search?q=pinarello+f12"

# abrir el navegador
browser = webdriver.Firefox()
# abrir la pagina
browser.get(uri)
# esperar hasta que se hayan renderizado los elementos que nos interesan (timeout=30s)
WebDriverWait(browser, 30).until(EC.presence_of_all_elements_located((By.CLASS_NAME, "rg_i.Q4LuWd.tx8vtf"))))
# obtener el código HTML
html = browser.page_source
# cerrar el navegador
browser.close()
```

6

Caso de Estudio: Descargar imágenes de Google Images

Se pretende descargarse automáticamente las imágenes devueltas en una búsqueda en Google Images y desarrollar un programa que las almacene en una carpeta local.



¿La página descargada coincide con la renderizada?

- En Firefox, lanzar una búsqueda en Google Images del término “pinarello F12”
<https://www.google.com/search?q=pinarello+f12&tbm=isch>
- Revisa el código HTML de la página descargada utilizando la opción del navegador “View page source”:

<https://www.google.com/search?q=pinarello+f12&tbm=isch>

[illegible]

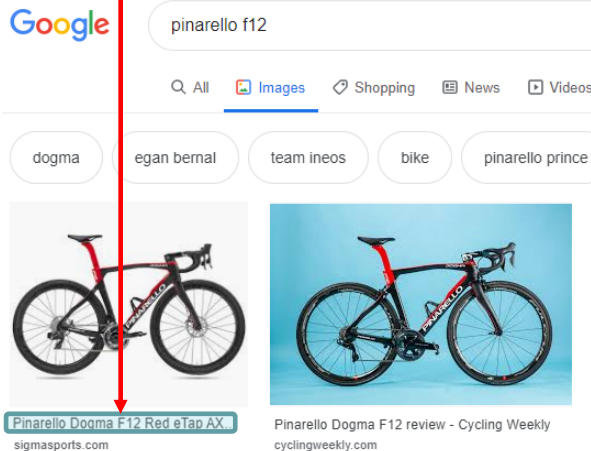
- El código contiene una pequeña estructura HTML básica y un gran cantidad de códigos JavaScript.
- Si se revisa el código con atención, se puede deducir que la función **_setImgSrc()** incluye las imágenes en el DOM del HTML.
- Por tanto, al no existir elementos `` con imágenes directamente enlazadas en la página descargada, no se puede utilizar la biblioteca requests para alimentar BeautifulSoup.

Analizar la estructura de la página HTML para saber dónde están las imágenes

- Por lo tanto, se deberá utilizar *selenium+geckodriver* para renderizar la hoja y poder parsear el html a DOM.
- Utilizando el bookmarklet “View Source Chart” buscaremos las imágenes dentro del código HTML

¿Dónde están las imágenes? Partimos de una cadena de texto que nos puede ayudar a ubicar una imagen, como por ejemplo, “F12 Red eTap”.

Buscamos esa cadena en el código HTML y analizamos la estructura HTML que la incluye. Las imágenes se encuentran en los elementos `` con un atributo "class" que tiene como valor `"rg_i Q4LuWd"`

[illegible]

Analizando en detalle la estructura de los elementos "img" que contienen los registros de resultados, se observa que hay dos tipos: a) los que tienen la imagen embebida en base64 y b) los que la enlazan

```
<div jaction="187Auh s9tbf,MW7te fL5lbf:RROdf s370ud,K3m4tCNbMv:03D1cIcyH7bi" data-ved="2aHUKetWgYozimeLnAIXD1eAKHVFSAmoQ4kDegQIARBa" data-toc="1" data-vid="LLK9Rz-9Y8EPM" jname="H9iUc" data-rn="31" class="lv-r PWICo MSM1f6 BuoofT jcntrl="S1436" jsmodel="Uz2pbF s84qc" jsdata="jOOpri:LLK9Rz-9Y8EPM-40" style="width: 272px; height: 222px;" data-tbuid="LLK9Rz-9Y8EPM" data-cv="0" data-cb="0" data-cl="0" data-cr="0" data-tv="273" data-cw="1200" data-cw="800">
  <a class="vWxWr isilb Eiy mM5pbf" jname="xTFXnd" jsaction="click:j9AeBi" data-nav="1" tabindex="0" style="height: 180px;">
    <div class="bRMDJf isilr" jname="bRMDJf isilr" data-ved="2aHUKetWgYozimeLnAIXD1eAKHVFSAmoQ4kDegQIARBa" data-toc="1" data-vid="LLK9Rz-9Y8EPM" jname="H9iUc" data-rn="31" class="lv-r PWICo MSM1f6 BuoofT jcntrl="S1436" jsmodel="Uz2pbF s84qc" jsdata="jOOpri:LLK9Rz-9Y8EPM-40" style="width: 272px; height: 222px;" data-tbuid="LLK9Rz-9Y8EPM" data-cv="0" data-cb="0" data-cl="0" data-cr="0" data-tv="273" data-cw="1200" data-cw="800">
      <div>
        <a data-arc="https://encrypted-tbn0.gstatic.com/images?q=tbn%3AAND9GQs012YldjeeC_Sg002QN25R772z6RvF1qk12Gc15dMFTT" data-rlt="w" jname="Q4LuWd" alt="Image result for pinarelo f12">
          </div>
        <div class="c7jWc">
          </div>
        <div class="h312td Rtl6" jname="bOER1">
          <div class="O1vY7">
            <span class="">
              1200 × 800
            </span>
          </div>
          <div class="PLlEc" jsaction="click: gfs2Re">
            </div>
          </div>
        </div>
      </div>
    <div class="vFACy kQAp" data-ved="2aHUKetWgYozimeLnAIXD1eAKHVFSAmoQ4kDegQIARBa" jname="uy6ald" rel="noopener" target="_blank" href="https://roadbikeaction.com/first-ride-pinarelo-dogma-f12-video/" jsaction="Focus:kvVbVb; mouseDown:kvVbVb; touchStart:kvVbVb">
      <div class="sM44c INtege">
        <div class="WGVvNb">
          First Ride: Pinarelo Dogma F12 - Video ...
        </div>
      </div>
    </div>
  </div>
</div>
```

[illegible]

Imagen *inline* en atributo "src"

6.4

Parsear el código HTML con BeautifulSoup

Para navegar por el árbol DOM del documento HTML utilizamos la librería *BeautifulSoup*:

<https://www.crummy.com/software/BeautifulSoup/bs4/doc/>

```
from bs4 import BeautifulSoup

# instanciar un parser para html y cargar en memoria el DOM del html
# "soup" es una ref. al elemento raíz del DOM
document = BeautifulSoup(html, 'html.parser')
# buscar en el DOM todos aquellos elementos cuyo atributo "class" valga "rg_i Q4LuWd tx8vtf"
img_results = document.find_all('img', {'class': 'rg_i Q4LuWd'})
for idx, each in enumerate(img_results):
    src = ""
    if each.has_attr('src'):
        src = each['src']
    else:
        src = each['data-src']
print(str(idx) + " " + src)
```

Almacenar imágenes

Realizamos peticiones HTTP para descargar las imágenes enlazadas y decodificamos las imágenes *inline* de base64 a binario.

```
import base64

# instanciar un parser para html y cargar en memoria el DOM del html
# "soup" es una ref. al elemento raíz del DOM
document = BeautifulSoup(html, 'html.parser')
# buscar en el DOM todos aquellos elementos cuyo atributo "class" valga "rg_i Q4LuWd"
img_results = document.find_all('img', {'class': 'rg_i Q4LuWd'})
for idx, each in enumerate(img_results):
    src = ""
    if each.has_attr('src'):
        src = each['src']
    else:
        src = each['data-src']
    print(str(idx) + " " + src)

img = None
if src.find("data:image") != -1:
    # data:[<mime type>][;charset=<charset>][;base64],<encoded data>
    img = base64.b64decode(src.replace("data:image/jpeg;base64,", ""))
else:
    res = requests.get(src)
    img = res.content

file = open("./img/" + str(idx) + ".jpeg", "wb")
file.write(img)
file.close()
```