

# Web Scraping

---

## Introducción, HTML

Mucha información está simplemente accesible en Internet. En este pequeño script vemos que es fácil acceder a ella, recopilarla, buscarla, procesarla.

El método más usado por los humanos para acceder a esa información de Internet es navegando mediante un navegador. El *Web Scraping* es la extracción de información que hay en sitios web, mediante programas de software

El contenido de las páginas web está escrito en archivos HTML. Veamos uno como ejemplo. En él, la información se organiza jerárquicamente: los distintos bloques de información están contenidos unos en otros, y se marcan con etiquetas. Mejor que cualquier explicación es verlo examinando el código fuente de cualquier página web. Empezamos con un archivo html muy pequeño:

```
In [1]: ▶ pequenno_codigo_html = """
<html>
  <head>
    <title>
      A web page
    </title>
  </head>
  <body>
    <p id="author">Joel Grus</p>
    <p id="subject">Data Science</p>
  </body>
</html>
"""

print(pequenno_codigo_html)
```

```
<html>
  <head>
    <title>
      A web page
    </title>
  </head>
  <body>
    <p id="author">Joel Grus</p>
    <p id="subject">Data Science</p>
  </body>
</html>
```

## BeautifulSoup

Empezamos por ver cómo se puede manejar un fragmento de texto html cualquiera, tanto si viene de Internet o si lo tenemos en nuestro equipo.

Para manejar la estructura de un archivo html, vamos a usar la librería BeautifulSoup. Esto requiere instalar el paquete bs4. Veamos las operaciones básicas.

```
In [2]: from bs4 import BeautifulSoup

soup_pequenno_codigo = BeautifulSoup(pequenno_codigo_html, "lxml")
primer_parrafo = soup_pequenno_codigo.find('p')

print(primer_parrafo.text)
print("-----")
print(primer_parrafo.text.split())
print(primer_parrafo["id"])
todos_los_parrafos = soup_pequenno_codigo.find_all('p')
print(todos_los_parrafos)
print(todos_los_parrafos[0].text)
```

Joel Grus

-----

['Joel', 'Grus']

author

[<p id="author">Joel Grus</p>, <p id="subject">Data Science</p>]

Joel Grus

Vamos ahora con una página web real:

<http://antares.sip.ucm.es/cpareja/>

Exáminala. Verás lo siguiente:



Cada navegador ofrece su propio modo de examinar el código fuente en html. En el mío, pulsando el botón derecho se puede seleccionar la opción *Ver el código fuente de la página*, y tenemos lo siguiente:



```
In [3]: ▶ import requests
requests.packages.urllib3.disable_warnings() # Prueba sin esta opción

mi_url = "https://antares.sip.ucm.es/cpareja/"
mi_codigo_html = requests.get(mi_url, verify=False).text # Prueba sin
print(mi_codigo_html[:250])
print(".....")
print(mi_codigo_html[-250:])
```

```
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Strict//EN" "http://www.w
3.org/TR/xhtml1/DTD/xhtml1-strict.dtd">

<!--
  Design by Free CSS Templates
  http://www.freecsstemplates.org (http://www.freecsstemplates.org)
  Released for free under a Creative Commons Attribution 2.5 L
  .....
  ga.src = ('https:' == document.location.protocol ? 'https://ssl' : 'h
  ttp://www') + '.google-analytics.com/ga.js';
    var s = document.getElementsByTagName('script')[0]; s.parentNode.in
  sertBefore(ga, s);
  })();

</script>
</body>

</html>
```

Y ahora vamos a usar funciones de la librería para extraer distintos tipos de elementos de dicha página.

```
In [4]: ▶ mi_pag_web = BeautifulSoup(mi_codigo_html, "lxml")
print(mi_pag_web.find('head'))
print("-----")
print(mi_pag_web.find('p'))
```

```
<head>
<meta content="" name="keywords"/>
<meta content="" name="description"/>
<meta content="text/html; charset=utf-8" http-equiv="content-type"/>
<title>CristÃ³bal Pareja Flores, pÃ¡gina web</title>
<link href="/style.css" media="screen" rel="stylesheet" type="text/cs
s"/>
</head>
-----
<p class="meta">
<span class="date">Curso 2013-14</span><span class="posted">Actualizad
o: 2014-abril-06</span>
</p>
```

```
In [5]: ▶ print(len(mi_pag_web.find_all('p')))
print("-----")
print(mi_pag_web.find_all('p')[8])
print("-----")
print(mi_pag_web.find_all('p', {'class': 'meta'}))
print("-----")
anclas_o_hiperenlaces = mi_pag_web.find_all('a')
print(len(anclas_o_hiperenlaces))
print(anclas_o_hiperenlaces[3])
print(anclas_o_hiperenlaces[3].get("href"))
print(anclas_o_hiperenlaces[3].get_text())

22
-----
<p class="meta">
<span class="date">2015</span><span class="posted">Actualizado: 2015-ab
ril-6</span>
</p>
-----
[<p class="meta">
<span class="date">Curso 2013-14</span><span class="posted">Actualizad
o: 2014-abril-06</span>
</p>, <p class="meta">
<span class="date">Hasta 2014</span><span class="posted">Actualizado: 2
014-dic-23</span>
</p>, <p class="meta">
<span class="date">2015</span><span class="posted">Actualizado: 2015-ab
ril-6</span>
</p>, <p class="meta">
<span class="date">Hasta 2014</span><span class="posted">Actualizado: 2
014-abril-6</span>
</p>]
-----
46
<a href="http://eprints.ucm.es/8705/1/CUADERNO_DE_TRABAJO_8.pdf">manual
</a>
http://eprints.ucm.es/8705/1/CUADERNO_DE_TRABAJO_8.pdf (http://eprints.
ucm.es/8705/1/CUADERNO_DE_TRABAJO_8.pdf)
manual
```

Vemos que los contenidos están estructurados en distintos contenedores identificados por distintas etiquetas en el documento `html`. Para acceder a estos contenidos, se han de usar funciones de la librería, de la manera que se ha mostrado, y con ellos se pueden extraer distintos tipos de elementos de la página y posteriormente procesarlos.

## Referencias:

Hay multitud de referencias a las funciones y métodos disponibles en esta librería, y hay multitud de ejemplos en Internet. Damos únicamente un par de dichas referencias.

1. Del libro "Data Science from Scratch", de Joel Grus, capítulo 9 ("Getting Data"), el apartado "Scraping the Web".
2. La página oficial de documentación de esta librería:  
<https://pypi.org/project/beautifulsoup4/>

