



Cristian, Abdel, Kevin

Pràctica 8.2: Web Scraping (XPath)

Lliuraments

Els resultats d'aquesta part de la pràctica s'hauran d'entregar en format PDF i l'entrega pot ser a través de GIT* o el moodle.

* S'ha d'entregar l'enllaç del GIT al moodle.

Guió

Amb l'ajuda de l'inspector d'elements del navegador, investiga com està formatada la pàgina <https://scrapepark.org/> . Aquesta pàgina està preparada per fer *web scraping*, de manera que les rutes per arribar als diferents elements no són trivials.

Exercici 1

Per començar, clona el repositori de GIT que es troba en aquesta ubicació i executa el codi Python per veure quin resultat dona.

https://github.com/pauitc/practica8_2

```
austria@PC002:~$ git clone https://github.com/pauitc/practica8_2
S'està clonant a «practica8_2»...
remote: Enumerating objects: 12, done.
remote: Counting objects: 100% (12/12), done.
remote: Compressing objects: 100% (5/5), done.
remote: Total 12 (delta 3), reused 12 (delta 3), pack-reused 0
S'estan rebent objectes: 100% (12/12), fet.
S'estan resolent les diferències: 100% (3/3), fet.
austria@PC002:~$ cd practica8_2
```

```
austria@PC002:~/practica8_2$ python3 web_scraping.py
<title>ScrapePark.org</title>
```

Exercici 2

- a. Executa les següents rutes XPath i observa el resultat que dona cada una. A continuació, explica les diferències que hi ha entre cada resultat i raona per què produeixen resultats diferents.

- i. node() vs text()

Ruta 1: `//div[@class='attribution']/p/node()`

'node()' retorna tots els nodes fills dels elements <p> que tenen una classe 'attribution'.

Ruta 2: `//div[@class='attribution']/p/text()`

'text()' retorna el contingut de text dels elements <p> que tenen una classe 'attribution'.

- ii. Barra simple vs barra doble

Ruta 1: `//ul[@class='navbar-nav']/li/a/text()`

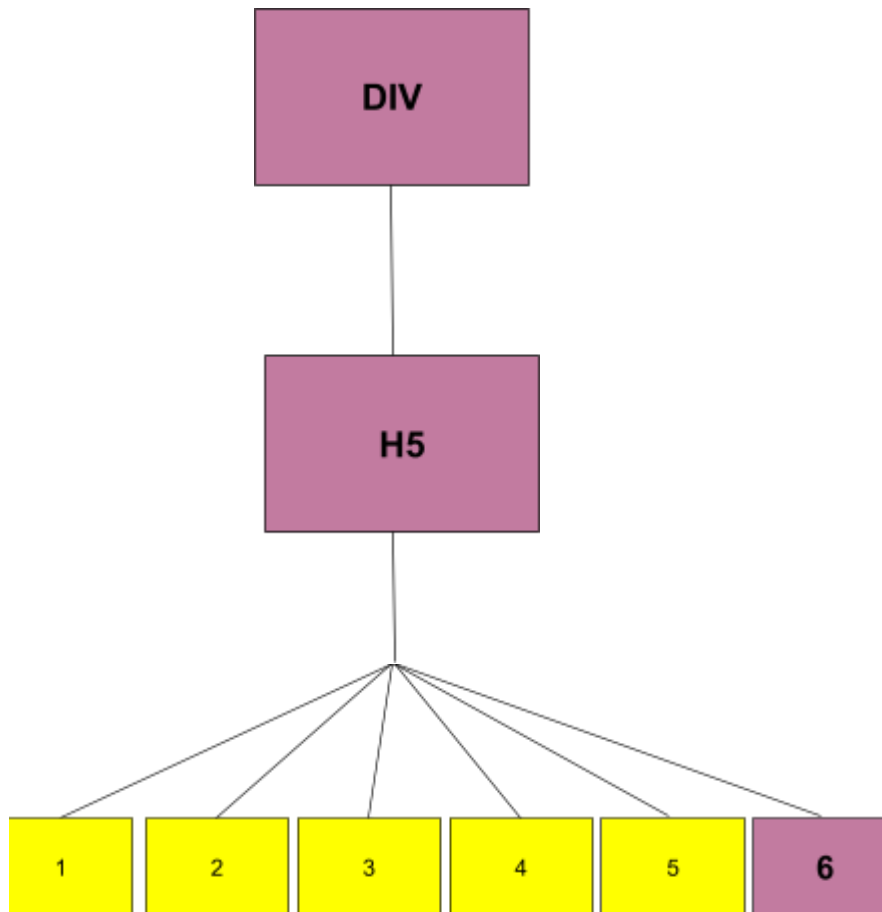
Busca tots els elements <a> que estan directament dins d'un element que és fill directe d'un element amb la classe 'navbar-nav', i retorna el seu contingut de text.

Ruta 2: `//ul[@class='navbar-nav']//li/a/text()`

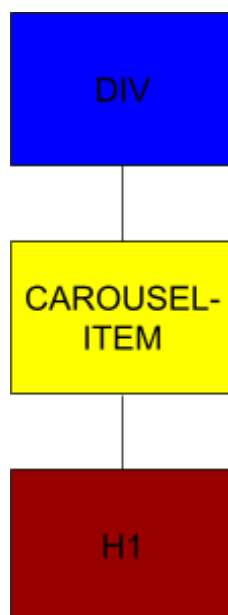
Busca tots els elements <a> que estan dins d'un element que pot ser fill directe o net dels elements amb la classe 'navbar-nav', i retorna el seu contingut de text.

- b. Representa, en forma d'arbre l'estructura HTML que resulta d'avaluar la següent ruta XPath (pots ignorar els salts de línia i espais).

- i. `(//div/h5)[6]`



ii. `//div[@class='carousel-item'][1]//h1`



Exercici 3

Descobreix la ruta XPath per arribar a cada un dels elements que es demana tenint en compte només la informació que es proporciona a l'enunciat.

- c. Troba la ruta que arriba al **correu** de contacte que es troba al **<footer>** de la pàgina.

Comença la ruta a l'etiqueta <html>

`/html`

sales@mail.com

`/html/body/footer/div/div/div/div/div/p[3]/span/text()`

- d. Troba la ruta que arriba a l'**atribut src** de la següent imatge (n'hi ha una al **<footer>**, i una al **<header>**, pots escollir):



images/logo.svg

`//header//a/img/@src`

- e. Troba la ruta fins a l'**atribut src** de les imatges amb **alt="Customer"**.

images/client-one.png

images/client-two.png

images/client-three.png

`//section[5]//div[@class='carousel-inner']//div[@class='img-box-inner']//img/@src`

- f. Troba la ruta fins a l'**adreça** de la pàgina web **"Fake Street 123"**. Fes que l'adreça XPath parteixi la següent ubicació:

`//div[@class='information-f']/p[1]/strong/text()`

Fake Street 123

`//div[@class='information-f']/p[1]/strong/text()/../..//span/text()`

- g. Troba la ruta que arriba fins al **<h5>** del “**New Skateboard 12**”. [**Pista**: busca la utilitat de la funció *normalize-space()*].

```
<h5>                                <span>New Skateboard</span> 12
</h5>
```

```
//section[3]//div[@class='row']/div[12]//div[@class='detail-box']/h5
/node()
```

- h. Partint de la ruta de l'apartat anterior, Troba la ruta que arriba fins al **preu** (text) del “**New Skateboard 12**”.

12

```
normalize-space(//section[3]//div[@class='row']/div[12]//div[@class=
'detail-box']/h5/text())
```

Exercici 4

Canvia la ruta a <https://scrapepark.org/table.html> . Amb l'ajuda del navegador, comprova què hi ha dins d'aquesta pàgina i troba la ruta XPath dels següents elements.

- i. Troba la ruta XPath a tots els **preus** dels **elements de color 'Blue'**. El resultat ha de ser el següent:

Blue
\$64
\$70
\$80
\$85

```
//tr[td[contains(text(),'Blue')]]/td/text()
```

- j. Troba la ruta que imprimeix **els preus del longboard** que es troben a la 4a columna de la taula **pintats en vermell**.

Longboard
\$80
\$85
\$90
\$62
\$150

```
//tr/td[4]/text()
```

- k. **Indica el nom i color** de l'article que **val \$110**. Comença l'expressió de la següent manera: [**pista**: hauràs de fer servir l'operador "|"]

```
//td[text()=' $110 ']
```

Skate
Special

```
//th[contains(text(),'Skate')]/text() | //tr[td[contains(text(), '$110')]]/td[1]/text()
```

- I. Troba la ruta a **tots els preus** dels objectes “Purple” **excepte el preu** que està pintat en vermell.

```
<td>Purple</td>  
<td class="text-center">$55</td>  
<td class="text-center">$60</td>  
<td class="text-center">$72</td>
```

```
//tr[td[contains(text(),'Purple')]]/td[not(contains(@style,'color:red'))]
```