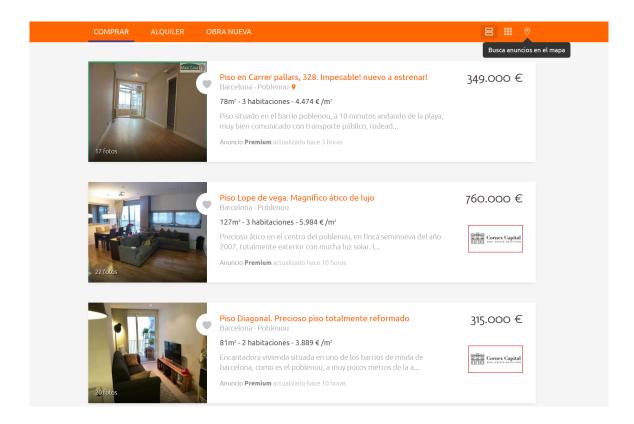# Tips for Web Scraping in Python with `bs4`

José María Lago Alonso

I've seen in the seminars that the less intuitive part of web-scraping via `bs4` appears when we try to scrap in the very deep of a nested structure, like in the following example:



Then the extra tools to scrap this data are:

```python
from urllib.request import urlopen
from bs4 import BeautifulSoup, NavigableString, Comment

url = "http://www.habitaclia.com/comprar-vivienda-en-barcelona-barrio_poblenou/provincia_ba
      rcelona-barcelones-area_6-sant_marti/listainmuebles.htm"

##Make the soup
page = urlopen(url) ## is up to your connection the URLError: <urlopen error [WinError
      10060]
soup = BeautifulSoup(page,"lxml")

casts = soup.find_all('ul', attrs={'class': 'enlista'})

cast = casts[0]
lista = []
for a in cast.find_all('a'):
    for child in a.children:
        if isinstance(child, NavigableString) and not isinstance(child, Comment) and
          str(child).strip() != "":
                lista.append('{}'.format(str(child).strip()))
```

```
18
19   newlist = []
20   for i in range(len(lista)):
21       if "si baja" in lista[i] or "a favoritos" in lista[i] or "Ver mapa" in lista[i]:
22           ""
23       else:
24           newlist.append(lista[i])
```

With this we capture the names of the adds in habitaclia.

The underlying idea of this example is to see the new objects `NavigableString` and also the `Comment`, so finally we have attributes like `children` that helps us to scrap in a nested way.