

# Chapitre V Web Scraping

*MI - IA / GL*

Ilyas Bambrik

# Table des matières



<b>Introduction</b>	3
<b>I - Web scraping HTML</b>	4
<b>II - BeautifulSoup</b>	8

# Introduction



Dans cette section nous nous intéressons à la collection de données avec deux modules particuliers. Premièrement, le module *requests* permettant de transmettre des requêtes HTTP et de recevoir les données renvoyées par le serveur. Par la suite, le module *BeautifulSoup* est utilisé pour extraire les informations pertinentes.



# Web scraping HTML



En science des données, nous pouvons réaliser de nombreux travaux avec le bon DataSet. Une fois que nous disposons de données intéressantes, nous pouvons utiliser Pandas ou Matplotlib pour analyser ou visualiser les tendances. Mais comment obtenir ces données en premier lieu ?

S'ils nous sont fournies dans un fichier csv ou json structuré / semi-structuré, le reste du travail est simple! Par contre, la plupart du temps, nous devons collecter les données nous-mêmes.

Souvent, vous trouverez le site Web idéal contenant toutes les données dont vous avez besoin. C'est là que *BeautifulSoup* est utile pour gratter (scrap) le code HTML. Si nous trouvons les données que nous souhaitons analyser en ligne, nous pouvons utiliser *BeautifulSoup* pour les transformer en une structure que nous pouvons comprendre. Cette bibliothèque Python nous permet de récupérer facilement et rapidement des informations d'un site Web et de les mettre dans un DataFrame.

Lorsque vous utiliser le Web scraping sur des sites Web, il est nécessaire de suivre certaines directives afin de respecter les sites et les conditions d'utilisation du site Web. Lisez la déclaration sur l'utilisation légale des données. Habituellement, les données que vous récupérez ne doivent pas être utilisées à des fins commerciales.

Ne spammez pas le site Web avec un tonne de requêtes. Un grand nombre de requêtes peut détruire un site Web qui n'est pas préparé à cette quantité de trafic. En règle générale, une bonne pratique est d'effectuer une requête sur une page Web chaque 30 secondes.

Par ailleurs, cette opération dépend de la structure du site web cible. Si la mise en page du site change, vous devrez modifier votre code de scraping pour être conforme à la nouvelle structure du site.

Python dispose d'une bibliothèque *requests* qui facilite grandement l'obtention de contenu. Tout ce que nous avons à faire est d'importer la bibliothèque, puis d'introduire l'URL du contenu que nous voulons obtenir. Le code suivant envoie une requête GET pour la ressource située dans la page <https://www.licbplus.com.dz/> :

```
1 import requests
2
3 webpage = requests.get('https://www.licbplus.com.dz/')
4 print(webpage.text)
```

Le résultat de l'exécution de ce programme est le code HTML retourné par le serveur web. Dans le contexte du Data Science, la collection des données périodique à partir d'un site distant serait intéressante pour créer un DataSet.

Idéalement le résultat serait structuré sous forme Json ou XML. La méthode *json()*, de l'objet retourné comme réponse, permet de transformer le résultat reçu en objet dictionnaire (qui a une syntaxe et structure identique à l'objet JSON). En outre, l'attribut *status\_code* peut être utilisé pour accéder au code réponse HTTP:

```
1 import requests
2
3 webpage = requests.get('https://googlechromelabs.github.io/chrome-for-testing/last-known-good-versions-with-downloads.json')
```

```
4 print(webpage.status_code)
5 print(webpage.json())
```

```
1 200
```

```
2 {'timestamp': '2024-07-24T11:08:54.116Z', 'channels': {'Stable': {'channel':
'Stable', 'version': '127.0.6533.72', 'revision': '1313161', 'downloads': {'chrome':
[{'platform': 'linux64', 'url': 'https://storage.googleapis.com/chrome-for-testing-
public/127.0.6533.72/linux64/chrome-linux64.zip'}, {'platform': 'mac-arm64', 'url':
'https://storage.googleapis.com/chrome-for-testing-public/127.0.6533.72/mac-arm64
/chrome-mac-arm64.zip'}, {'platform': 'mac-x64', 'url': 'https://storage.googleapis.
com/chrome-for-testing-public/127.0.6533.72/mac-x64/chrome-mac-x64.zip'},
{'platform': 'win32', 'url': 'https://storage.googleapis.com/chrome-for-testing-public
/127.0.6533.72/win32/chrome-win32.zip'}, {'platform': 'win64', 'url':
'https://storage.googleapis.com/chrome-for-testing-public/127.0.6533.72/win64/chrome-
win64.zip'}], 'chromedriver': [{'platform': 'linux64', 'url': 'https://storage.
googleapis.com/chrome-for-testing-public/127.0.6533.72/linux64/chromedriver-linux64.
zip'}, {'platform': 'mac-arm64', 'url': 'https://storage.googleapis.com/chrome-for-
testing-public/127.0.6533.72/mac-arm64/chromedriver-mac-arm64.zip'}, {'platform':
'mac-x64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/127.
0.6533.72/mac-x64/chromedriver-mac-x64.zip'}, {'platform': 'win32', 'url':
'https://storage.googleapis.com/chrome-for-testing-public/127.0.6533.72/win32
/chromedriver-win32.zip'}, {'platform': 'win64', 'url': 'https://storage.googleapis.
com/chrome-for-testing-public/127.0.6533.72/win64/chromedriver-win64.zip'}], 'chrome-
headless-shell': [{'platform': 'linux64', 'url': 'https://storage.googleapis.com
/chrome-for-testing-public/127.0.6533.72/linux64/chrome-headless-shell-linux64.zip'},
{'platform': 'mac-arm64', 'url': 'https://storage.googleapis.com/chrome-for-testing-
public/127.0.6533.72/mac-arm64/chrome-headless-shell-mac-arm64.zip'}, {'platform':
'mac-x64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/127.
0.6533.72/mac-x64/chrome-headless-shell-mac-x64.zip'}, {'platform': 'win32', 'url':
'https://storage.googleapis.com/chrome-for-testing-public/127.0.6533.72/win32/chrome-
headless-shell-win32.zip'}, {'platform': 'win64', 'url': 'https://storage.googleapis.
com/chrome-for-testing-public/127.0.6533.72/win64/chrome-headless-shell-win64.
zip'}]}], 'Beta': {'channel': 'Beta', 'version': '127.0.6533.57', 'revision':
'1313161', 'downloads': {'chrome': [{'platform': 'linux64', 'url': 'https://storage.
googleapis.com/chrome-for-testing-public/127.0.6533.57/linux64/chrome-linux64.zip'},
{'platform': 'mac-arm64', 'url': 'https://storage.googleapis.com/chrome-for-testing-
public/127.0.6533.57/mac-arm64/chrome-mac-arm64.zip'}, {'platform': 'mac-x64', 'url':
'https://storage.googleapis.com/chrome-for-testing-public/127.0.6533.57/mac-x64
/chrome-mac-x64.zip'}, {'platform': 'win32', 'url': 'https://storage.googleapis.com
/chrome-for-testing-public/127.0.6533.57/win32/chrome-win32.zip'}, {'platform':
'win64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/127.
0.6533.57/win64/chrome-win64.zip'}], 'chromedriver': [{'platform': 'linux64', 'url':
'https://storage.googleapis.com/chrome-for-testing-public/127.0.6533.57/linux64
/chromedriver-linux64.zip'}, {'platform': 'mac-arm64', 'url': 'https://storage.
googleapis.com/chrome-for-testing-public/127.0.6533.57/mac-arm64/chromedriver-mac-
arm64.zip'}, {'platform': 'mac-x64', 'url': 'https://storage.googleapis.com/chrome-
for-testing-public/127.0.6533.57/mac-x64/chromedriver-mac-x64.zip'}, {'platform':
'win32', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/127.
0.6533.57/win32/chromedriver-win32.zip'}, {'platform': 'win64', 'url':
'https://storage.googleapis.com/chrome-for-testing-public/127.0.6533.57/win64
/chromedriver-win64.zip'}], 'chrome-headless-shell': [{'platform': 'linux64', 'url':
'https://storage.googleapis.com/chrome-for-testing-public/127.0.6533.57/linux64
/chrome-headless-shell-linux64.zip'}, {'platform': 'mac-arm64', 'url':
'https://storage.googleapis.com/chrome-for-testing-public/127.0.6533.57/mac-arm64
/chrome-headless-shell-mac-arm64.zip'}, {'platform': 'mac-x64', 'url':
'https://storage.googleapis.com/chrome-for-testing-public/127.0.6533.57/mac-x64
/chrome-headless-shell-mac-x64.zip'}, {'platform': 'win32', 'url': 'https://storage.
googleapis.com/chrome-for-testing-public/127.0.6533.57/win32/chrome-headless-shell-
win32.zip'}, {'platform': 'win64', 'url': 'https://storage.googleapis.com/chrome-for-
testing-public/127.0.6533.57/win64/chrome-headless-shell-win64.zip'}]}], 'Dev':
{'channel': 'Dev', 'version': '128.0.6601.2', 'revision': '1328598', 'downloads':
{'chrome': [{'platform': 'linux64', 'url': 'https://storage.googleapis.com/chrome-for-
testing-public/128.0.6601.2/linux64/chrome-linux64.zip'}, {'platform': 'mac-arm64',
'url': 'https://storage.googleapis.com/chrome-for-testing-public/128.0.6601.2/mac-
arm64/chrome-mac-arm64.zip'}, {'platform': 'mac-x64', 'url': 'https://storage.
googleapis.com/chrome-for-testing-public/128.0.6601.2/mac-x64/chrome-mac-x64.zip'},
{'platform': 'win32', 'url': 'https://storage.googleapis.com/chrome-for-testing-public
/128.0.6601.2/win32/chrome-win32.zip'}, {'platform': 'win64', 'url': 'https://storage.
googleapis.com/chrome-for-testing-public/128.0.6601.2/win64/chrome-win64.zip'}],
'chromedriver': [{'platform': 'linux64', 'url': 'https://storage.googleapis.com
/chrome-for-testing-public/128.0.6601.2/linux64/chromedriver-linux64.zip'},
{'platform': 'mac-arm64', 'url': 'https://storage.googleapis.com/chrome-for-testing-
public/128.0.6601.2/mac-arm64/chromedriver-mac-arm64.zip'}, {'platform': 'mac-x64',
'url': 'https://storage.googleapis.com/chrome-for-testing-public/128.0.6601.2/mac-x64
/chromedriver-mac-x64.zip'}, {'platform': 'win32', 'url': 'https://storage.googleapis.
com/chrome-for-testing-public/128.0.6601.2/win32/chromedriver-win32.zip'},
{'platform': 'win64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public
/128.0.6601.2/win64/chromedriver-win64.zip'}], 'chrome-headless-shell': [{'platform':
'linux64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/128.
0.6601.2/linux64/chrome-headless-shell-linux64.zip'}, {'platform': 'mac-arm64',
'url': 'https://storage.googleapis.com/chrome-for-testing-public/128.0.6601.2/mac-
arm64/chrome-headless-shell-mac-arm64.zip'}, {'platform': 'mac-x64', 'url':
'https://storage.googleapis.com/chrome-for-testing-public/128.0.6601.2/mac-x64/chrome-
headless-shell-mac-x64.zip'}, {'platform': 'win32', 'url': 'https://storage.
googleapis.com/chrome-for-testing-public/128.0.6601.2/win32/chrome-headless-shell-
googleapis.com/chrome-for-testing-public/128.0.6601.2/win32/chrome-headless-shell-
```

```
win32.zip'}, {'platform': 'win64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/128.0.6601.2/win64/chrome-headless-shell-win64.zip'}}}, 'Canary': {'channel': 'Canary', 'version': '129.0.6614.0', 'revision': '1332009', 'downloads': {'chrome': [{'platform': 'linux64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/linux64/chrome-linux64.zip'}, {'platform': 'mac-arm64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/mac-arm64/chrome-mac-arm64.zip'}, {'platform': 'mac-x64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/mac-x64/chrome-mac-x64.zip'}, {'platform': 'win32', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/win32/chrome-win32.zip'}, {'platform': 'win64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/win64/chrome-win64.zip'}]}, 'chromedriver': [{'platform': 'linux64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/linux64/chromedriver-linux64.zip'}, {'platform': 'mac-arm64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/mac-arm64/chromedriver-mac-arm64.zip'}, {'platform': 'mac-x64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/mac-x64/chromedriver-mac-x64.zip'}, {'platform': 'win32', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/win32/chromedriver-win32.zip'}, {'platform': 'win64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/win64/chromedriver-win64.zip'}]}, 'chrome-headless-shell': [{'platform': 'linux64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/linux64/chrome-headless-shell-linux64.zip'}, {'platform': 'mac-arm64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/mac-arm64/chrome-headless-shell-mac-arm64.zip'}, {'platform': 'mac-x64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/mac-x64/chrome-headless-shell-mac-x64.zip'}, {'platform': 'win32', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/win32/chrome-headless-shell-win32.zip'}, {'platform': 'win64', 'url': 'https://storage.googleapis.com/chrome-for-testing-public/129.0.6614.0/win64/chrome-headless-shell-win64.zip'}]}}
```

Comme déclaré initialement, il est préférable de transmettre des requêtes sur des intervalles régulières. Pour ceci, la fonction `sleep()` du module `time` est utilisée. Par exemple, le code suivant illustre comment transmettre une requête à l'URL précédant une fois par 24 heures. La fonction `time.sleep()` prend le temps d'attente en secondes :

```
1 import requests
2 import time
3 while True:
4     webpage = requests.get('https://googlechromelabs.github.io/chrome-for-testing
5     /last-known-good-versions-with-downloads.json')
6     print(webpage.status_code)
7     print(webpage.json())
8     time.sleep(24*60*60)
```

Dans les deux premiers exemples, nous avons utilisé `.text` et `.json()` de l'objet *réponse* afin d'accéder / transformer le contenu reçu. En cas où le résultat reçu n'est pas un contenu textuel (image, fichier exe, etc), l'attribut `.content` est utilisé afin d'accéder à la suite d'octets brute. La valeur détenue par cette attribut est de type *bytes*.

De même, l'attribut `.headers` du résultat contient les entêtes renvoyés par le serveur (Content-type, Last-Modified, etc) sous forme de dictionnaire :

```
1 import requests
2 webpage = requests.get('https://www.licbplus.com.dz/')
3 print(webpage.headers)
4 print(webpage.headers['Content-Type'])

1 {'Connection': 'Keep-Alive', 'Keep-Alive': 'timeout=5, max=100', 'content-type':
  'text/html; charset=UTF-8', 'link': '<https://www.licbplus.com.dz/wp-json/>; rel="
  https://api.w.org/", <https://www.licbplus.com.dz/wp-json/wp/v2/pages/36>; rel="
  alternate"; type="application/json", <https://www.licbplus.com.dz/>; rel="shortlink",
  'etag': '"600708-1721822742;br"', 'x-litespeed-cache': 'hit', 'content-encoding':
  'br', 'vary': 'Accept-Encoding', 'content-length': '36888', 'date': 'Wed, 24 Jul 2024
  16:51:41 GMT', 'server': 'DZ-HTTP'}
2 text/html; charset=UTF-8
3 {'Connection': 'Keep-Alive', 'Keep-Alive': 'timeout=5, max=100', 'content-type':
  'text/html; charset=UTF-8', 'link': '<https://www.licbplus.com.dz/wp-json/>; rel="
  https://api.w.org/", <https://www.licbplus.com.dz/wp-json/wp/v2/pages/36>; rel="
  alternate"; type="application/json", <https://www.licbplus.com.dz/>; rel="shortlink",
  'etag': '"600708-1721822742;br"', 'x-litespeed-cache': 'hit', 'content-encoding':
  'br', 'vary': 'Accept-Encoding', 'content-length': '36888', 'date': 'Wed, 24 Jul 2024
  16:51:41 GMT', 'server': 'DZ-HTTP'}
```

```
4 text/html; charset=UTF-8
```

On peut aussi spécifier les entêtes de la requête transmis au serveur par le paramètre optionnel *headers* :

```
1 import requests
2 import time
3
4 webpage = requests.get('https://googlechromelabs.github.io/chrome-for-testing
5 /last-known-good-versions-with-downloads.json',
6 headers={'Accept': 'application/json'})
7 print(webpage.status_code)
8 print(webpage.json())
```

Le module *requests* propose une fonction pour chaque méthode HTTP (GET, HEAD, POST, PUT, DELETE, OPTIONS). Généralement, s'il est nécessaire de transmettre des données au serveur (comme la clé d'authentification ou le contenu d'un formulaire) pour recevoir le résultat, ces derniers sont exprimés dans le corps de la requête au lieu de l'url et la méthode POST est utilisée. Le paramètre optionnel *data* est utilisé pour exprimer les données à transmettre sous forme de dictionnaire, liste de tuples, d'octets ou un objet fichier :

```
1 import requests
2 objDict= {'somekey': 'somevalue'}
3 resultat = requests.post('https://www.exemple.API.com/', data= objDict, headers= {
4 "Cookie": "=XKJHQSf21924JSAQF"})
5 print(resultat .text)
```

### Remarque : API

Une page web qui retourne un résultat sous forme JSON/XML (contenu semi-structurée) est appelée *API (Application Programming Interface)*. Par exemple l'API suivante vous renvoie des information de localisation géographique sur une adresse IP donnée (161.185.160.93 dans l'exemple)

<https://ipinfo.io/161.185.160.93/geo>

<https://ipinfo.io/geo> (sans le segment contenant l'adresse IP), l'API renvoie des informations sur la machine qui a générer la requête.

# BeautifulSoup



Le traitement des données semi-structurées JSON ou XML est facile généralement. Par contre, il est plus difficile de récupérer les données à partir d'un contenu HTML. Par ailleurs, si la structure de la page change, il faudra adapter le code afin d'accéder aux informations ciblées. *BeautifulSoup* est une bibliothèque qui facilite le traitement d'un document HTML pour récupérer le contenu cible. Supposons que nous sommes intéressés par les prix des PC bureau proposés par licbplus.com. Pour commencer, il faudra inspecter le code HTML/structure de la page. Les prix sont présentés dans un élément *ul* définit avec la classe *product\_list\_widget*.

```

138 av ecome-menu-wapper"><ul id="menu-primary-menu" class="clone-main-menu ecome-clone-mobile-menu ecome-nav main-menu"><li id="menu-item-1
139 rated_products"><h2 class="widgettitle">Les meilleurs articles<span class="arrow"></span></h2><ul class="product_list_widget"><li>
140
141 -RTX-1-300x300.jpg" class="attachment-woocommerce_thumbnail size-woocommerce_thumbnail" alt="" loading="lazy" srcset="https://www.licbpl
142 class="rating">5.00</strong> sur 5</span></div><span class="review">( 1 review(s) )</span></div>
143 mbol">DA</span></span></li><li>
144
145 1-300x300.jpg" class="attachment-woocommerce_thumbnail size-woocommerce_thumbnail" alt="" loading="lazy" srcset="https://www.licbplus.c
146 class="rating">5.00</strong> sur 5</span></div><span class="review">( 1 review(s) )</span></div>
147 ol">DA</span></span></li><li>
148
149 -300x300.jpg" class="attachment-woocommerce_thumbnail size-woocommerce_thumbnail" alt="" loading="lazy" srcset="https://www.licbplus.com
150 class="rating">5.00</strong> sur 5</span></div><span class="review">( 1 review(s) )</span></div>
151 ol">DA</span></span></li><li>
152
153 super-300x300.jpg" class="attachment-woocommerce_thumbnail size-woocommerce_thumbnail" alt="" loading="lazy" srcset="https://www.licbpli
154 class="rating">5.00</strong> sur 5</span></div><span class="review">( 1 review(s) )</span></div>

```

*BeautifulSoup* utilise les concepts de navigation DOM (Document Object Model) HTML qui permettent la sélection des éléments selon l'identifiant, classe css et attributs. Pour commencer, il est nécessaire d'importer le module *BeautifulSoup* et puis de parser le contenu HTML. Le type de paramètre à parser doit être spécifier dans le deuxième argument qui permet aussi de parser d'autres types de contenu. Par la suite, le résultat retourné est similaire à l'objet *document* en Javascript qui permet la sélection d'élément dans le document. Le code suivant utilise la méthode *.select()* pour retrouver les éléments *li* contenus les éléments *ul* possédants la classe *product\_list\_widget*.

Tout de même, les éléments retournés sont toujours des éléments HTML. Pour accéder au texte contenu dans chaque élément, il suffit d'utiliser l'attribut *.text* :

```

1 import requests
2 from bs4 import BeautifulSoup
3 webpage = requests.get('https://www.licbplus.com.dz/product-category/pc/pc-aures/'
4 )
5 bs = BeautifulSoup(webpage.content, "html.parser")
6 result=bs.select("ul.product_list_widget>li")
7 print(result)
8 print([li.text for li in result])
9

```

```

1 LICB+ I101 RTX
2 Note 5.00 sur 5 ( 1 review(s) )
3 110,900DA
4
5 HAVIT HV-H2106d
6 Note 5.00 sur 5 ( 1 review(s) )
7 1,790DA

```



```

8
9 HAVIT MS1022
10 Note 5.00 sur 5( 1 review(s) )
11 1,690DA
12
13 LICB+ A36 SUPER
14 Note 5.00 sur 5( 1 review(s) )
15 108,900DA
16
17 HAVIT H2026D
18 Note 5.00 sur 5( 1 review(s) )
19 2,790DA

```

Il est possible d'accéder aux attributs de chaque élément retourné par la méthode `.select()` avec la méthode `.get()`. Le code suivant sélectionne les images contenues dans la page, puis pour chaque image `.get('src')` renvoie le chemin vers l'image cible.

```

1 import requests
2 from bs4 import BeautifulSoup
3 webpage = requests.get('https://www.licbplus.com.dz/product-category/pc/pc-aures/'
4 )
5 bs = BeautifulSoup(webpage.content, "html.parser")
6 images=bs.select("img")
7 print(*map(lambda ele:ele.get("src"),images))

```

La méthode `.select()` admet comme paramètre n'importe quel sélecteur `css` valide comme `document.querySelectorAll()` en Javascript. Le code suivant sectionne les éléments ayant un attribut `data-testid="atom-body-paragraph"` qui sont des descendants quelconques d'un élément ayant l'attribut `data-testid="molecule-live-comment"`.

```

1 reponse=requests.get("https://www.eurosport.fr/football/transferts/2024-2025
2 /liveevent.shtml")
3 document = BeautifulSoup(reponse.content, "html.parser")
4 live_feed=document.select('[data-testid="molecule-live-comment"] [data-testid="
5 atom-body-paragraph"]')
6 [news.text.strip() for news in live_feed]

```

```

1 ['Simons, retour au point de départ ?',
2 "Nouveau rebondissement concernant le cas Xavi Simons. Joueur sous contrat avec
3 le Paris Saint-Germain, l'international néerlandais est intéressé par un nouveau
4 prêt, après celui de la saison passée à Leipzig. Selon les derniers échos relayés par
5 Fabrizio Romano, le Batave serait finalement proche... d'un retour à Leipzig ! Le
6 club allemand ferait le forcing pour obtenir un nouveau prêt. Une possibilité qui
7 pourrait convaincre Paris.",
8 '',
9 'Yvon Mvogo vers Bournemouth\xa0?',
10 'Sous contrat jusqu'en juin 2025 avec Lorient, le portier Yvon Mvogo pourrait
11 rebondir du côté de Bournemouth en Premier League selon l'Equipe. Les deux clubs
12 partagent le même propriétaire ce qui pourrait faciliter la transaction. Mvogo
13 bénéficie d'un bon de sortie accordé par son président Loïc Féry.',
14 '',
15 'Deuxième offre du Bayern pour Désiré Doué',
16 "Le Bayern Munich vient de proposer environ 55 millions d'euros bonus compris à
17 Rennes pour recruter Doué selon l'Equipe. Une offre qui pourrait satisfaire les
18 Rennais qui en réclame un peu près 60 millions. Le PSG devrait revenir à la charge et
19 proposer une nouvelle offre à la hauteur du club allemand.",
20 '',
21 'Morgan Guilavogui en prêt à Sankt Pauli',
22 'L'attaquant Morgan Guilavogui va être prêté avec option d'achat par Lens à
23 Sankt Pauli, tout juste promu en Bundesliga selon l'Equipe.',
24 '',
25 'Brentford va vendre Ivan Toney',

```

14 'Selon The Athletic, Brentford va vendre Ivan Toney lors de ce mercato, à un an  
de la fin de son contrat. Les Bees préfèrent recevoir une indemnité de transfert cet  
été plutôt que de le voir partir gratuitement l'année prochaine.',  
15 '',

16 'Assignon pisté par Aston Villa',

17 'Le latéral droit du Stade Rennais, Lorenz Assignon, prêté la saison dernière à  
Burnely, intéresse le club anglais d'Aston Villa selon Foot Mercato. Rennes attend au  
moins 10 millions d'euros pour lâcher son joueur de 24 ans, aussi ciblé par l'AS  
Rome.',  
18 '',

19 'Sarr d'accord avec Crystal Palace',

20 'Selon The Athletic, Ismaila Sarr s'est mis d'accord avec Crystal Palace pour  
rejoindre le club. Un contrat de quatre ou cinq ans l'attend à Londres. Palace doit  
maintenant négocier avec Marseille.',  
21 '',

22 'Le PSG explore la piste Sancho',

23 'Selon Sky Sports, le PSG n'a pas encore proposé d'offre à Manchester United  
pour s'attacher les services de Jadon Sancho. Le club de la capitale étudie pour le  
moment le dossier. La Juventus, Dortmund ainsi qu'un autre club de Premier League  
sont également sur le coup.',  
24 '',

25 'Séville veut Ansu Fati',

26 'Selon les informations de Sport, relayées par Foot Mercato, le FC Séville fait  
le forcing pour recruter Ansu Fati. Le Barça veut le vendre et Séville a pour le  
moment proposé une offre en dessous des attentes du club catalan. Les négociations se  
poursuivent.',  
27 '',

28 'Le PSG pousse pour Doué et cible Williams',

29 'Désireux de recruter le milieu offensif de Rennes Désiré Doué, Nasser Al-  
Khelaifi s'investit personnellement dans le dossier selon The Athletic. Pour le  
moment, il n'y a aucun accord entre le club parisien et le Stade Rennais.  
Parallèlement, le PSG a noué les premiers contacts avec Bilbao pour le recrutement de  
Nico Williams. Sous contrat jusqu'en 2027 avec le club espagnol, l'ailier dispose  
d'une clause libératoire d'environ 60 millions d'euros.',  
30 '',

31 'Abdoulaye Ndiaye va signer à Brest',

32 'À la recherche d'un renfort en défense pour remplacer Lilian Brassier qui est  
parti à l'OM, Brest va recruter le défenseur central de Troyes Abdoulaye Ndiaye.  
Selon Foot Mercato, la transaction est un prêt avec option d'achat de 9 millions  
d'euros.',  
33 '',

34 'Marseille bientôt d'accord avec Nketiah',

35 'À la recherche d'un attaquant pour pallier le départ de Pierre-Emerick  
Aubameyang, l'OM est proche d'un accord contractuel avec Eddie Nketiah selon  
l'Equipe. Il faut encore convaincre son club d'Arsenal avec qui il est sous contrat  
jusqu'en 2027.',  
36 '',

37 'Adams signe au Torino',

38 'Libre depuis l'expiration de son contrat avec Southampton, Che Adams vient de  
s'engager librement avec le Torino. C'est la première expérience à l'étranger pour  
l'Écossais.',  
39 '',

40 'L'Inter Milan lorgne Zeze',

41 'Le champion d'Italie en titre l'Inter Milan a approché Nantes pour s'attacher  
les services de Nathan Zeze, jeune défenseur central de 19 ans selon l'Equipe. Sous  
contrat jusqu'en 2028 avec les Canaris, Zeze est considéré comme intransférable par  
la famille Kita.',  
42 '',

43 'Chelsea sur Filip Jorgensen',

44 'Chelsea a approché le portier de Villarreal Filip Jorgensen pour un éventuel  
transfert selon The Athletic. Le gardien de but danois est également ciblé par  
l'Olympique de Marseille.',  
45 '',

46 'Savic va signer à Trabzonspor',

47 'Stefan Savic va quitter librement l'Atlético de Madrid pour s'engager avec  
Trabzonspor selon Relevo. Il va signer un contrat jusqu'en juin 2027 en Turquie.',  
48 '',

49 'De Bruyne va rester à Manchester City',

```

50 'Pep Guardiola a confirmé que le milieu belge Kevin de Bruyne va rester à
Manchester City cette saison malgré les rumeurs d'un potentiel transfert en Arabie
saoudite.',
51 '',
52 'Wan-Bissaka sur le départ',
53 'Selon les informations de Sky Sport Germany, Manchester United pousse pour le
départ d'Aaron Wan-Bissaka. United demande environ 10 millions d'euros pour se
séparer de l'Anglais. Les Red Devils ciblent Noussair Mazraoui pour le remplacer.',
54 '',
55 'Bonjour à tous et bienvenue dans ce nouveau live pour suivre les dernières
rumeurs du mercato estival',
56 'Le mercato en direct - Sancho, Varane, Todibo, PSG, OM, Real Madrid : Suivez la
journée des transferts en direct sur Eurosport']

```

Alternativement à `.select()`, la recherche peut être faite avec la méthode `.find()` et `.find_all()`. Il est aussi possible de limiter le nombre de résultats retournés par `select()`, `find` et `find_all`, avec l'argument optionnel `limit`.

Le paramètre optionnel string de `find_all` peut être utilisé avec une fonction qui filtre les résultats retournés. Le code suivant recherche les éléments `li` et puis, avec la fonction passée dans l'argument string, seulement les éléments contenant la chaîne "HDD" sont retournés :

```

1 import requests
2 from bs4 import BeautifulSoup
3 webpage = requests.get('https://www.licbplus.com/dz/product-category/pc/pc-aures/'
)
4 bs = BeautifulSoup(webpage.content, "html.parser")
5 elements=bs.find_all("li",string=lambda text:"HDD" in text if text!=None else
False )
6 print(elements)
7

```

```

1 [<li class="menu-item menu-item-type-custom menu-item-object-custom menu-item-
4268" id="menu-item-4268"><a class="ecome-menu-item-title" href="https://www.licbplus.
com/product-category/composants-pc/hdd-interne/" title="HDD interne">HDD interne</a><
/li>,
2 <li class="menu-item menu-item-type-custom menu-item-object-custom menu-item-
12739" id="menu-item-12739"><a class="ecome-menu-item-title" href="https://www.
licbplus.com/product-category/composants-pc/rack-hdd-interne/" title="Rack HDD
interne">Rack HDD interne</a></li>,
3 <li class="menu-item menu-item-type-custom menu-item-object-custom menu-item-
4269" id="menu-item-4269"><a class="ecome-menu-item-title" href="https://www.licbplus.
com/product-category/peripherique-pc/hdd-externe/" title="HDD externe">HDD externe<
/a></li>]

```