

Collecte Automatisée de Données Web

Séance 2 de cours/TD

IUT SD Dole

Objectifs:

- | | | |
|---|--|---|
| — Introduction à la collecte de données Web sans API... | | — |
|---|--|---|

Webscraping : Qu'est-ce que c'est ? _____[COURS]

Le webscraping désigne les techniques d'extraction du contenu des sites internet. C'est une pratique très utile pour toute personne souhaitant travailler sur des informations disponibles en ligne, mais pas nécessairement accessible par une API.

L'idée générale va être de charger le code source de la page web puis d'aller chercher dans celui-ci les informations souhaitées.

Commençons par un rappel rapide de la structure d'un site web : Un site Web est un ensemble de pages codées en HTML qui permet de décrire à la fois le contenu et la forme d'une page Web.

Sur une page web, vous trouverez toujours à coup sûr des éléments comme `<head>`, `<title>`, etc. Il s'agit des codes qui vous permettent de structurer le contenu d'une page HTML et qui s'appellent des balises. Citons, par exemple, les balises `<p>`, `<h1>`, `<h2>`, `<h3>`, `` ou ``. Le symbole `< >` est une balise : il sert à indiquer le début d'une partie. Le symbole `</ >` indique la fin de cette partie. La plupart des balises vont par paires, avec une balise ouvrante et une balise fermante (par exemple `<p>` et `</p>`).

Pour récupérer correctement les informations d'un site internet, il faut pouvoir comprendre sa structure et donc son code HTML. Les fonctions python qui servent au scraping sont principalement construites pour vous permettre de naviguer entre les balises.

Webscraping en python _____[COURS]

De manière assez similaire à l'utilisation d'API, faire du webscraping en python commence par faire la requête de la page web.

La différence principale est qu'à ce moment-la, plutôt que de récupérer des données pré-formatées dans un format json, xml ou csv, on récupère le code brut de la page web, qu'il va falloir explorer pour obtenir les données voulues.

Pour cela, on va se servir de la bibliothèque BeautifulSoup. (D'autres existent, notamment Selenium, qui vous permettra de faire du webscraping sur les parties dynamiques des pages, sur lequel on reviendra plus bas). Son rôle va être de structurer la donnée récupérée afin d'accéder facilement à son contenu en filtrant sur les balises par exemple. Prenons un cas particulier. Nous allons essayer de récupérer le tableau de données de la page wikipédia suivante : https://fr.wikipedia.org/wiki/Liste_des_pays_par_population

On commence par regarder le code source de la page pour repérer le tableau. Vous le trouverez en ligne 481 dans une balise `table`. En y regardant de plus près, on notera que chaque ligne est délimitée par une balise `tr`, et chaque case à l'intérieur de la ligne par une balise `td`

```
1 import requests
2 import pandas as pd
3 import bs4
4
5 #On recupere le contenu brut de la page
6 url_api = "https://fr.wikipedia.org/wiki/Liste_des_pays_par_population"
7 req = requests.get(url_api)
8
9 # On utilise beautifulsoup pour formater le code source
10 page = bs4.BeautifulSoup(req.text, "html.parser")
11
12 #On recupere le tableau
13 tab = page.find("table")
14
15 #Puis on traite chaque ligne
16 res=[]
17 for line in tab.find_all("tr"):
18     elems = line.find_all("td")
19     resline=[]
20     for elem in elems:
21         resline=resline +[elem.text.strip()]
22     res=res+[resline]
23
24 print(res)
```

la méthode `get_text` permet de récupérer le contenu textuel d'une balise en enlevant aussi toutes les balises internes.

Pour plus de détails sur les fonctionnalités de la bibliothèque BeautifulSoup, référez-vous à sa documentation : <https://www.crummy.com/software/BeautifulSoup/bs4/doc/>

Les Bonnes Pratiques de Webscraping [COURS]

Le webscraping est une activité qui peut causer un certain nombre de problèmes aux sites ciblés, ce qui en fait quelque chose d'assez mal vu en général. Pour faire du webscraping de manière polie et raisonnable, il y a quelques règles à suivre :

- Comme pour les API, assurez-vous de limiter vos chargements de pages au strict nécessaire.
- Insérez des pauses entre les différents appels pour ne pas surcharger les serveurs (En python, vous pourrez utiliser la méthode `sleep` de la bibliothèque `time`.
- Souvenez-vous qu'il est illégal de webscraper des données considérées comme personnelles, même si elles sont publiées sur un site.
- Respectez les souhaits des webmaster du site concernant ce qui est, ou non, webscrapable sur le site en question.

Pour ce dernier point, vous pourrez par exemple vérifier ce qui est indiqué sur le `robots.txt` du site web.

`robots.txt` est un fichier disponible normalement sur tout site web (à l'adresse `http://www.adressedusite.com/robots.txt`) et indiquant ce qu'il est possible de faire sur le site. Initialement créé pour les robots de moteur de recherche, il est aussi utilisé pour indiquer les limites du webscraping.

Par exemple, l'extrait du fichier robots suivant :

```
1 User-agent: *
2 Disallow: /mandataire/*
3 Disallow: /coach/*
4 Disallow: /notreequipe/*
5
6 User-agent: bingbot
7 Crawl-delay: 1
```

indique que les pages du site dans le sous-dossier mandataire sont interdites à tout robot (et donc à tout web-scraping), et que les robots de type bingbot doivent respecter un délai d'une seconde entre deux pages chargées.

Il est fortement conseillé de suivre ces recommandations, pour éviter d'être banni des sites en question.

Exercice 1 (,)

Afficher la liste des liens présents dans la page wikipedia proposée dans l'exemple.

□

Exercice 2 (,)

En allant récupérer des données sur wikipedia, produisez un fichier contenant la liste des rois de france avec leur date de règne.

□

Exercice 3 (,)

En allant récupérer des données sur des sites de votre choix, proposez un programme permettant de vérifier automatiquement si deux acteurs donnés ont déjà joué ensemble.

□

Exercice 4 (,)

En allant récupérer des données sur trois sites d'info de votre choix, proposez un programme permettant de générer un fichier csv affichant les titres et les liens vers les infos du jour.

□

Exercice 5 (,)

Ecrire un programme qui fonctionne comme un dictionnaire (l'utilisateur entre un mot et le programme affiche les définitions du mot concerné).

□

Webscraper des pages dynamiques _____[COURS]

Nous avons précédemment vu comment récupérer les données de pages statiques (où l'ensemble des données est présent dans le code html de la page à étudier). Nous allons maintenant voir comment récupérer les données pages dynamiques, contenant notamment du javascript. Nous allons pour cela utiliser la bibliothèque Selenium de python.

Cette bibliothèque permet de simuler une connection depuis un navigateur web, et de manipuler la page comme pourrait le faire un utilisateur. En particulier, il est possible pour le programme python de scroller pour charger plus de la page, de cliquer sur des boutons, ou de remplir des formulaires, notamment d'authentification.

Vous aurez besoin d'installer la bibliothèque selenium, mais aussi la bibliothèque webdriver_manager pour simplifier les choses.

Prenons un exemple :

```
1 # On commence par charger les bibliothèques
2
3 from selenium import webdriver
4 from selenium.webdriver.chrome.service import Service
5 from webdriver_manager.chrome import ChromeDriverManager
6 from selenium.webdriver.common.keys import Keys
7 from selenium.webdriver.common.by import By
8 from time import sleep
9
10
11 # On créer un simulateur de navigateur web, ici chrome
12 driver = webdriver.Chrome(service=Service(ChromeDriverManager().install
13   ()))
14
15 # On connecte le simulateur à google
16 driver.get("https://www.google.com/")
17
18
19 # On detecte le bouton "accepter les cookies"
20 cookiebutton = driver.find_element(by=By.ID, value="L2AGLb")
21
22 # on clique sur ce bouton
23 cookiebutton.click()
24
25 # On patiente quelques secondes que la page se charge entièrement
26 sleep(3)
27
28 # On detecte la barre de recherche de google
29 search = driver.find_element(by=By.NAME, value="q")
30
31
32 # on y entre le texte selenium, puis on appuie sur entree
33 search.send_keys("Selenium")
34 search.send_keys(Keys.ENTER)
35 sleep(3)
36
37 # on scrolle tout en bas de la page, pour charger la deuxième page de résultats
38 driver.execute_script("window.scrollTo(0, document.body.scrollHeight);")
39 sleep(3)
40
41 # on detecte tous les elements html de balise h3 ( qui contiennent les titres des
42   résultats de recherche google
43
44 # on affiche les textes de ces resultats.
45 for x in (links2):
46     print(x.text)
```

on a ainsi récupéré les résultats des deux premières pages google pour la recherche "selenium".

Exercice 6 (,)

Récupérez sur le site www.la-maronne-immobiliere.com l'ensemble des informations standard sur les biens immobiliers situés à moins de 50 km d'Aurillac. □