

# DATASCIENCE, LEARNING AND APPLICATIONS

DALAS - Data acquisition via web scrapping

20 janvier 2024

Laure Soulier - Nicolas Baskiotis

INSPIRED FROM OLIVIER SCHWANDER'S LECTURE  
AND INSEE'S "FORMATION WEB SCRAPPING"

## Web scrapping

- **Web Scrapping** : content extraction of websites
- **Web crawling** : search of new URLs on internet or on a specific website

But both can collaborate : crawler for gathering URLs and then, scrapping for extracting the content.

### Examples

- Price follow-up, temporal analysis
- Sentiment analysis (comment on products, feedbacks, ...)
- Potential clients (Pages Jaunes, Google Maps, ...)
- Meta-engines (lastminute.com, kiwi.com), gathering info from different websites

## Is it legal ?

- Not illegal, but the use of data might be regulated (GDPR, CNIL, intellectual property law, data gouvernance act, jurisprudence, licences, ...)
- Rules depend on the country you are and the data you scrap.
- Sell or forward scrapped data might be illegal

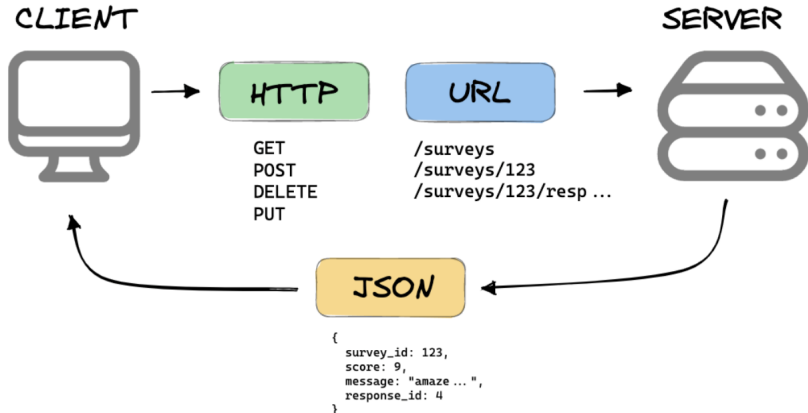
### L'article 323-3 du code pénal

"Le fait [...] d'extraire, de détenir, de reproduire [...] frauduleusement les données qu'il [le site] contient est puni de cinq ans d'emprisonnement et de 150 000e d'amende"

## Basic notions of the web

- Client-server architecture (request/answer)
- Web service : Application/software function allowing an inter-operable communication over the network between different platforms/tools. They use standard and open protocols.
- Most used architecture style : Representational State Transfer (REST). The URI is important (allows to identify resources and values, available operations : GET, POST, PUT, DELETE, ...)
- Others protocols/architecture styles SOAP and WSDL (more complexed and more structured)

# WHAT IS A REST API?



## Format

`http://example.com/chemin/vers/la/ressource?arg1=valeur1&arg2=valeur2`

- `http` : protocol
- `example.com` : server
- `path/to/the/resource` : identifier of the resource
- `?` : everything that follows is an arg
- `arg1=valeur1&arg2=valeur2` : args

## Requesting API Rest

- We use the HTTP protocol to call CRUD actions :

```
//Get all devices: HTTP GET
http://api.example.com/device-management/managed-devices
//Create new Device: HTTP POST
http://api.example.com/device-management/managed-devices
//Get device for given Id: HTTP GET
http://api.example.com/device-management/managed-devices/{id}
//Update device for given Id: HTTP PUT
http://api.example.com/device-management/managed-devices/{id}
//Delete device for given Id: HTTP DELETE
http://api.example.com/device-management/managed-devices/{id}
```

For more info : <https://restfulapi.net/resource-naming/>



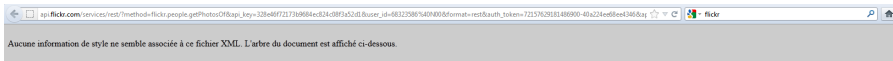
```

http://search.twitter.com/search.json?q=blue%20angels&rpp=5&include_entities=true&result_type=mixed - Mozilla Firefox
Fichier  Edition  Affichage  Historique  Marque-pages  Outils  ?
http://search.twitter.com/search.json... x  JSONView :: Modules pour Firefox  x  +
search.twitter.com/search.json?q=blue%20angels&rpp=5&include_entities=true&result_type=mixed  Google

from_user: "BreezyRawr",
from_user_id: 399249177,
from_user_id_str: "399249177",
from_user_name: "Tara ",
geo: null,
id: 178089017538134000,
id_str: "178089017538134016",
iso_language_code: "en",
- metadata: {
  result_type: "recent"
},
profile_image_url: http://a0.twimg.com/profile_images/1857194391/426058_115188978607231_100003483870163_61170_472193261_n_normal.jpg,
profile_image_url_https: https://si0.twimg.com/profile_images/1857194391/426058_115188978607231_100003483870163_61170_472193261_n_normal.jpg,
source: "&lt;a href=&quot;http://twitter.com/&quot;&gt;web&lt;/a&gt;",
text: "RT @AmyLovesCMB_Xx: Roses are red, Violets are blue, God sent us angels, Yes Breezy that's you. ♥",
to_user: null,
to_user_id: null,
to_user_id_str: null,
to_user_name: null
},
- {
  created_at: "Fri, 09 Mar 2012 12:01:18 +0000",
  - entities: {
    hashtags: [ ],
    urls: [ ],
    user_mentions: [ ]
  },
  from_user: "AmyLovesCMB_Xx",
  from_user_id: 63974669,
  from_user_id_str: "63974669",
  from_user_name: "amy robinson ",
  geo: null,
  id: 178088065171066880,
  id_str: "178088065171066880",

```

## API REST : Flickr



```

- <rsp stat="ok">
- <photos page="1" pages="4" perpage="20" total="70" has_next_page="1">
  <photo id="6818874820" owner="29893924@N08" secret="75a1264454" server="7199" farm="8" title="Feliz Dia de la Mujer." ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6857161799" owner="34829221@N00" secret="e142ceae339" server="7202" farm="8" title="What would you ask?" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6852635517" owner="8705967@N02" secret="725d8198cf" server="7152" farm="8" title="V de Venezuela, V de VOTA!!" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6689735273" owner="48664097@N02" secret="b1735cdd60" server="7021" farm="8" title="Llamada" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6642489581" owner="25257508@N03" secret="a4699ec6d4" server="7025" farm="8" title="Things I Love Thursday" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6641187815" owner="48739681@N05" secret="9c5144ca5a" server="7016" farm="8" title="~ Feliz Nit de Reis ~ Feliz Noche de Reyes ~" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6596007895" owner="34829221@N00" secret="9e16d2d366" server="7165" farm="8" title="2011 in Review" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6593866295" owner="38838924@N03" secret="f867e832d4" server="7153" farm="8" title="Mi 2011 en flickr" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6590116673" owner="8705967@N02" secret="23bd13d96a" server="7018" farm="8" title="Navidad" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6587696631" owner="8332135@N03" secret="63cc7bd96a" server="7142" farm="8" title="The sky is the limit..." ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6586167143" owner="63307805@N00" secret="cfaebd377" server="7007" farm="8" title="my year in pictures" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6581955027" owner="39307146@N08" secret="7b56e63043" server="7171" farm="8" title="My year in pictures" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6566465059" owner="23523125@N08" secret="4c73a0eb85" server="7147" farm="8" title="Feliz Navidad..." ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6564959719" owner="29219049@N06" secret="4200050be9d" server="7153" farm="8" title="Feliz Navidad!!!" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6560620287" owner="40654531@N03" secret="800a0970a" server="7001" farm="8" title="mi navidad, tu navidad? Explored Dec 23, 2011 #6 Muchas Gracias!" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6553379789" owner="49763065@N08" secret="20b628586" server="7002" farm="8" title="." ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6538048271" owner="37763325@N08" secret="66c29587d" server="7171" farm="8" title="Felices Fiestas! ~ Happy Holidays!" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6526003925" owner="66881369@N06" secret="e6bb093294" server="7020" farm="8" title="Family lunch" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6475299671" owner="43616712@N03" secret="aa811060f" server="7146" farm="8" title="lovely support" ispublic="1" isfriend="0" isfamily="0"/>
  <photo id="6358869783" owner="25257508@N03" secret="205188c38" server="6019" farm="7" title="Friday Favorites" ispublic="1" isfriend="0" isfamily="0"/>
</photos>
</rsp>

```

## Request and Answer format

### → HTTP Request

- Other element in the request : meta-data, encoding, authentication, cookies, session tracking, ...
- Answer format : CML, JSON, Binary, proprietary, ...

### → Authentication mechanisms

- Login + password sent every time : rare
- Login + password sent once then secret : token mechanism
- Secret in URL : `https://example.com/e1WnQ10c` (Doodle for instance)
- Two factor login : sms, email, ...

- **Requests - HTTP for Humans**

*Warning : Recreational use of other HTTP libraries may result in dangerous side-effects, including : security vulnerabilities, verbose code, reinventing the wheel, constantly reading documentation, depression, headaches, or even death.*

(source : <http://docs.python-requests.org/en/master/>)

- **HTTPX - A next-generation HTTP client for Python**

Same but looks better

## Requesting REST API in Python

→ Ideal case with web service, documentation provided, open formats, authorized access, ...

### Base

```
1 >>> import requests
2 >>> result = requests.get("http://www.bing.com/search?
    q=http")
3 >>> print(result.text)
```

### Json : Request does data loading

```
1 >>> import requests
2 >>> r = requests.get('https://api.github.com/events')
3 >>> r.json()
4 [{u'repository': {u'open_issues': 0, u'url': 'https://
    github.com/...
```

- If no documentation, we do reverse engineering : analysis of web page requests, study of mobile applications

But often...

→ Less-than-ideal case : no web service, incomprehensible format, incompatible terms of use, not the information you're looking for.

## What to do ?

- Data extraction from a web page
  - No guarantee of stability
  - Not necessarily authorized / general terms of use
  - Responsible behavior essential

## How ?

- Make HTTP requests
- Extract data from responses

### **woob : Web Out Of Browser <https://woob.tech>**

- Interfaces for other sites
- Mainly for banks, but not only
- Professional support

### **Zotero <https://www.zotero.org/>**

- Bibliographic database
- Collect metadata of scientific publications
- largely used by researchers

## Examples in production

### Scrapy <https://scrapy.org>

- Easy website browsing

### Playwright <https://playwright.dev/>

- Automates a webbrowser

### Web crawling

- websearch engine
- Archiving, tracking, ...



## The "Classic" web

- GET request on the URL
- HTML code as answer : the page is a document
- The document is a HTML tree : easy to capture/browse

## Libraries

- Requesting API : Requests / HTTPX
- Parsing HTML : BeautifulSoup

## Example of HTML code

```
1 <!DOCTYPE html>
2 <html>
3   <head>
4     <title>
5       Exemple de HTML
6     </title>
7   </head>
8   <body>
9     This is a sentence with an
10      <a href="cible.html">hyperlink</a>.
11     <p>
12 This is a paragraph without hyperlink </p>
13   </body>
14 </html>
```

## XPath

- Request language on XML tree
- Localization path : succession of steps to reach a target
- Localization step : information filtering and extraction

## CSS

- Set up for page display
- But is an instance of tree browsing

## Examples

### XPath

---

/	root node
/html/body	body node
//div	all div nodes
//h1[@class='title']	nodes of a given class
div	relative request
div[1]/p[2]	Second paragraph of first div
table[count(tbody/tr)>2]	Tables with more than two lines

---

### XPath

- div.class#id
- and other filters

## Python lib for syntax analysis (XML, HTML)

```
1 >>> import requests
2 >>> from bs4 import BeautifulSoup
3 >>> url = "https://fr.wikipedia.org/wiki/
    Uniform_Resource_Locator"
4 >>> page = requests.get(url)
5 >>> print(page.status_code)
6     200
7 >>> soup = BeautifulSoup(page.text, "html.parser")
8 >>> soup.title
9     <title>Uniform Resource Locator \u2014 Wikip\u00e9dia
    </title>
10 >>> soup.title.string
11     u'Uniform Resource Locator \u2014 Wikip\u00e9dia'
12 >>> soup.find("h2")
13     <h2>Sommaire</h2>
```

# In Python - BeautifulSoup

## Find a node

```
1 >>> soup.find("div")
2 >>> soup.find("div", class="offer")
3 >>> soup.find("div", id="firstHeading")
```

## Find all nodes

```
1 >>> soup.find_all("div")
```

## Parents, enfants

```
1 >>> find_parents() , .find_parent()
2 >>> .content[0]
3 >>> .next_siblings , .previous_siblings
```

## The "modern" web

- GET request on the URL
- A piece of information
- Javascript code : the page is a program
- The code executes other requests and build the page accordingly
- The content is on a virtual tree (DOM) : the tree is the result of a program execution
- More complex !

## Libraries

- Selenium (when javascript : allow to simulate a browser and to interact with javascript elements)
- Scrapy (more complete and more complex tool. Recommended for big projects)

## WebDriver / remote control interface

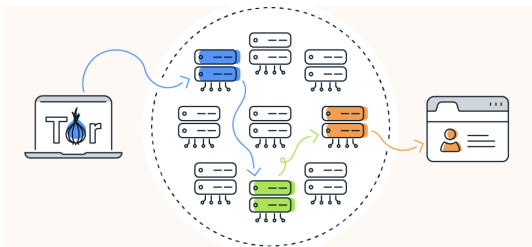
```
1 >>> driver = webdriver.Chrome()
2 >>> driver.get("https://www.selenium.dev/selenium/web/
  web-form.html")
3 >>> title = driver.title
4 >>> text_box = driver.find_element(by=By.NAME, value="
  my-text")
5 >>> submit_button = driver.find_element(by=By.
  CSS_SELECTOR, value="button")
```



```
1 >>> from playwright.sync_api import sync_playwright
2 >>> with sync_playwright() as p:
3     browser = p.chromium.launch()
4     page = browser.new_page() page.goto("http://
playwright.dev")
5     print(page.title()) browser.close()
```

## Scrapping on a private network

- TOR (<https://www.torproject.org/fr/download/> and <https://www.lobstr.io/fr/blog/comment-scrap-avec-des-tor-proxies-et-python>)  
Tor (The Onion Router) is a network that anonymizes web traffic to ensure 100% private web browsing. The Tor browser masks your IP address and browsing activity by redirecting web traffic via a series of routers known as nodes.



→ **This does not prevent you from respecting the law**  
Introduction, data science and web scrapping

## Protections against scraping

- Check that the user is a human
  - Captcha
  - Two-factor authentication
  - Solution.... A human....
- Statistics
  - Rate limit : freeze when the speed of interactions is too high
  - Machine learning to detect non-human behavior
  - Solution : look like a human
- Other technics
  - Frequently updating the page structure
  - IP filtering

## Conclusion

- Libraries
  - Complete navigators : Playwright, Selenium
  - Just request : Scrapy, Requests, HTTPS, BeautifulSoup
  - Private navigation : Tor
- Tools
  - Inspect the navigator
  - Tree browsing : XPath and CSS
- No formal method :
  - Understanding the website
  - Mimic human behavior
  - Responsible usage to limit legal/ethical issues

More info : <https://inseefrlab.github.io/formation-webscraping/>