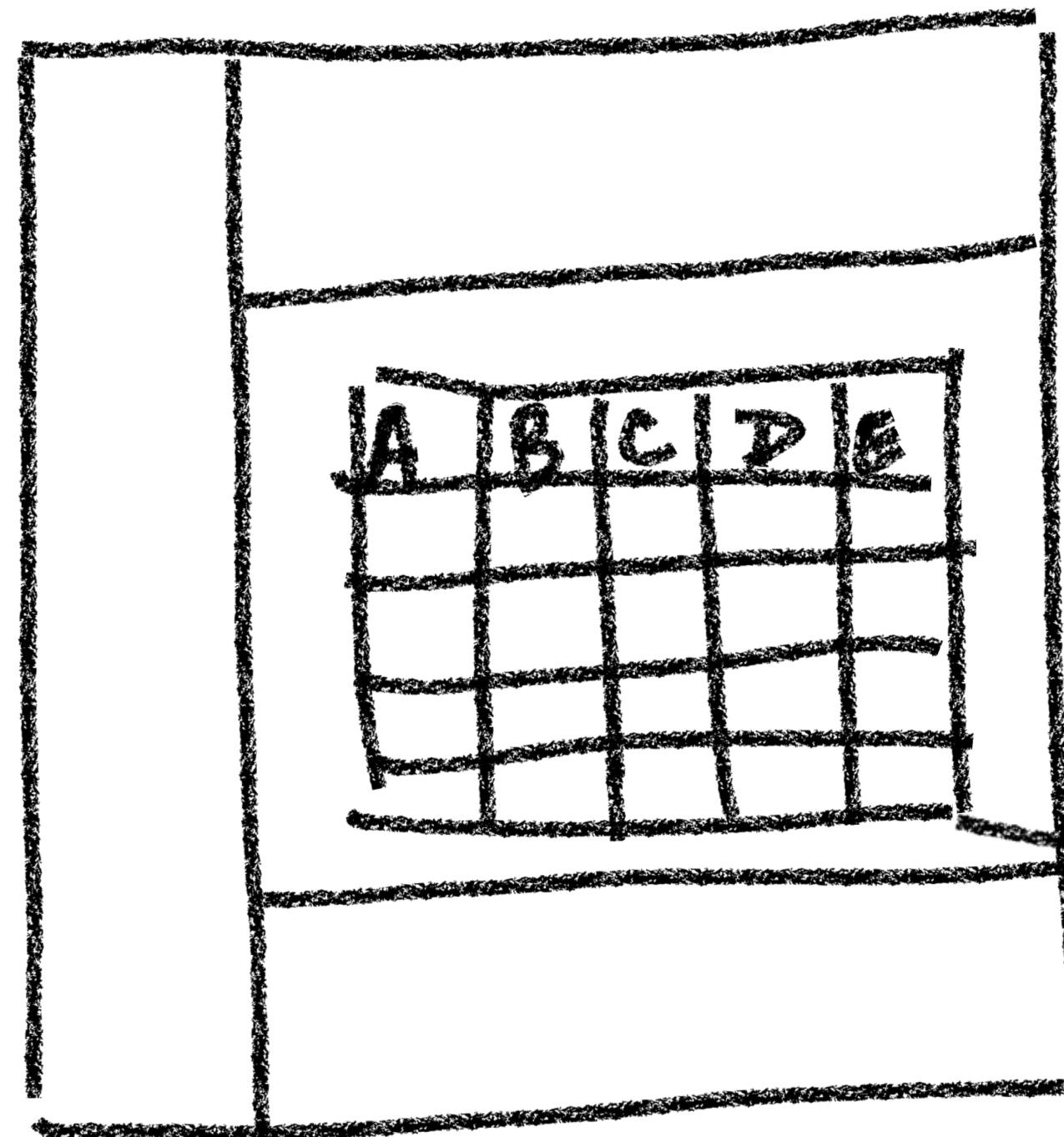


## Web Scraping

Aspirer les données d'une page

web de manière structurée.



`<html>`  
`<head>`  
`</head>`  
`<body>`  
`</body>`  
`</html>`



Data Frame

	A	B	C	D	E
0	1	2	3	4	5
1	6	7	8	9	10
2	11	12	13	14	15
i					
N					

Code HTML

## ① Aspirer le code source :

(= Récupérer le code source sous la forme  
d'un string)

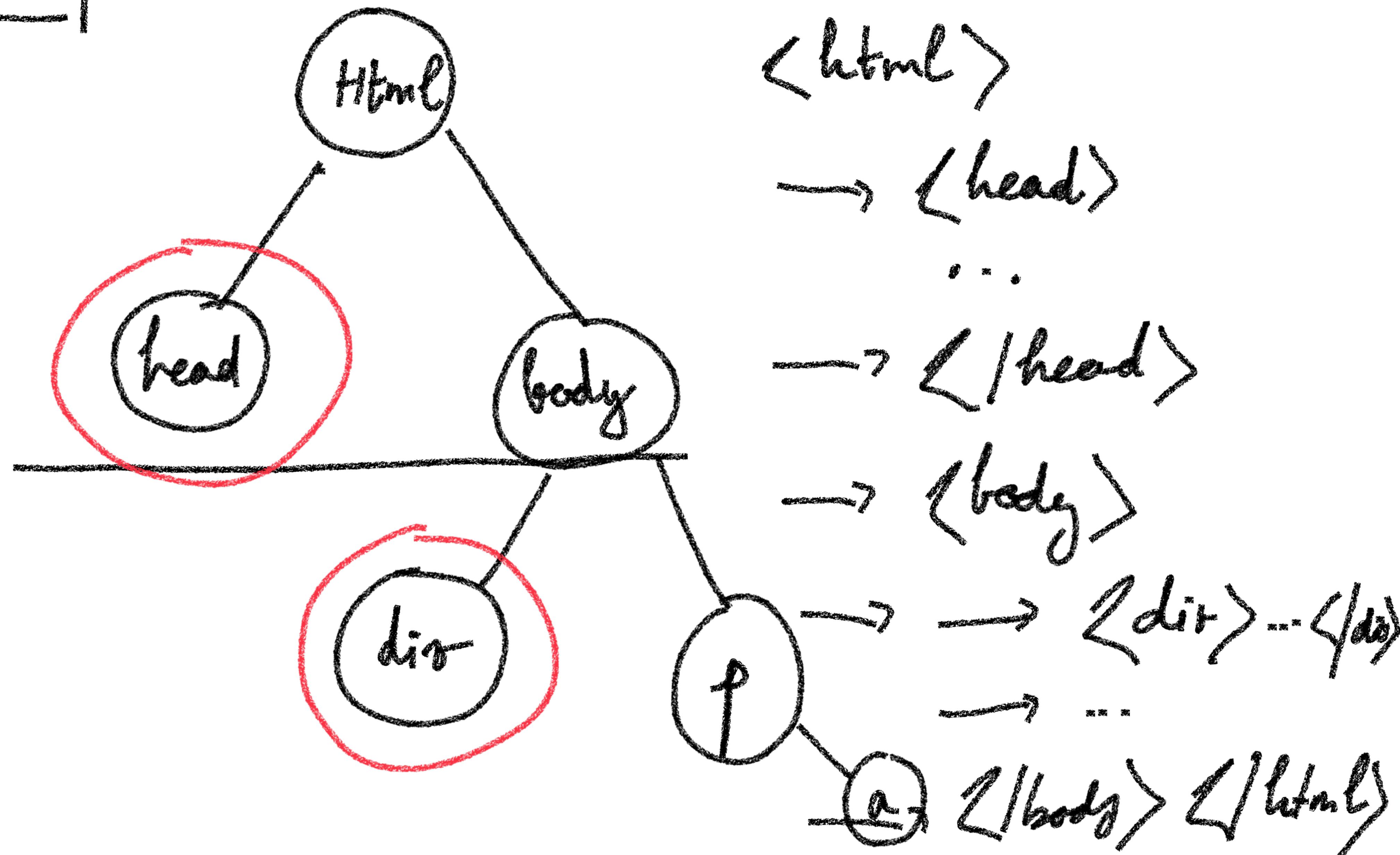
Requête GET de l'URL de la page

```
> res = requests.get("http://x.y.z/a/b/c")
```

```
> res.content
```

Dom HTML

Document - Object Model = ABLE



## ① Analyser / Parser le code HTML

### BeautifulSoup

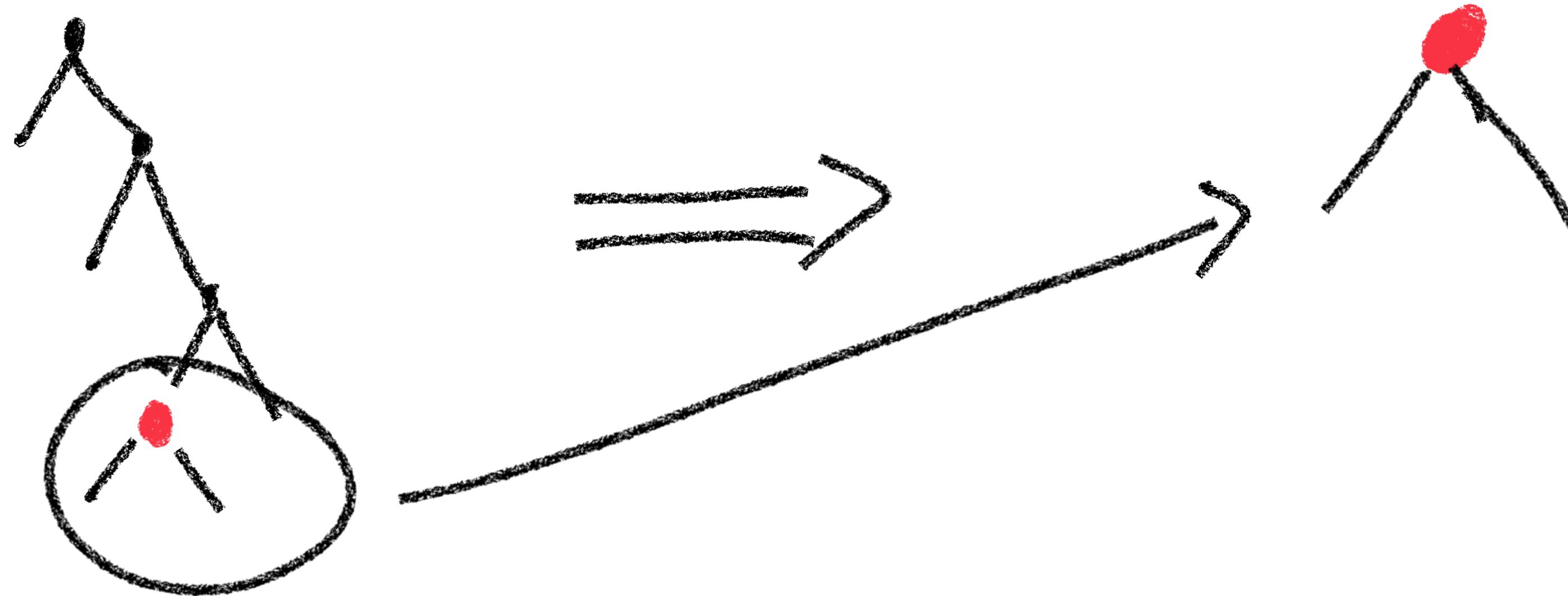
soup = BeautifulSoup(rs.content, "html")

#### ② Recherche ciblée d'informations:

soup.select\_one(...): renvoie la première balise q'il trouve

soup.select\_all(...): renvoie la liste de toutes les balises qui matchent

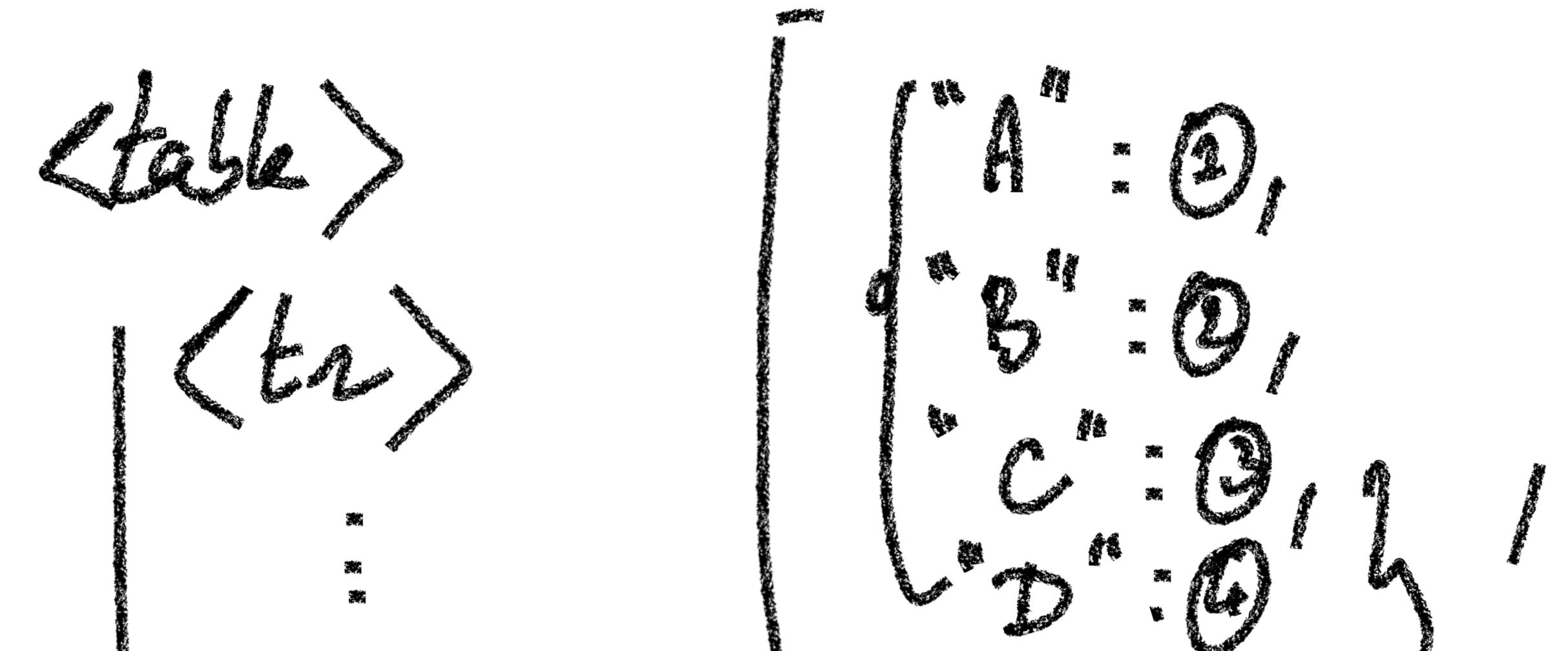
Soup.select \*  $\Rightarrow$  un objet de type  
BeautifulSoup  
G. select \*



dc

## Constructing the DataFrame:

A	B	C	D
①	②	③	④



pd.DataFrame(  
from\_dict  
(...))

## ① Limitations légales:

https://www.x.y/robots.txt

User-Agent: *
Allowed: /a/b
Disallowed: /c/d

User  
Agent  
= Navigator

## ① Limitations techniques:

### ②a) Limites de quota de requêtes:

Par ex : une seule requête / 10 s.

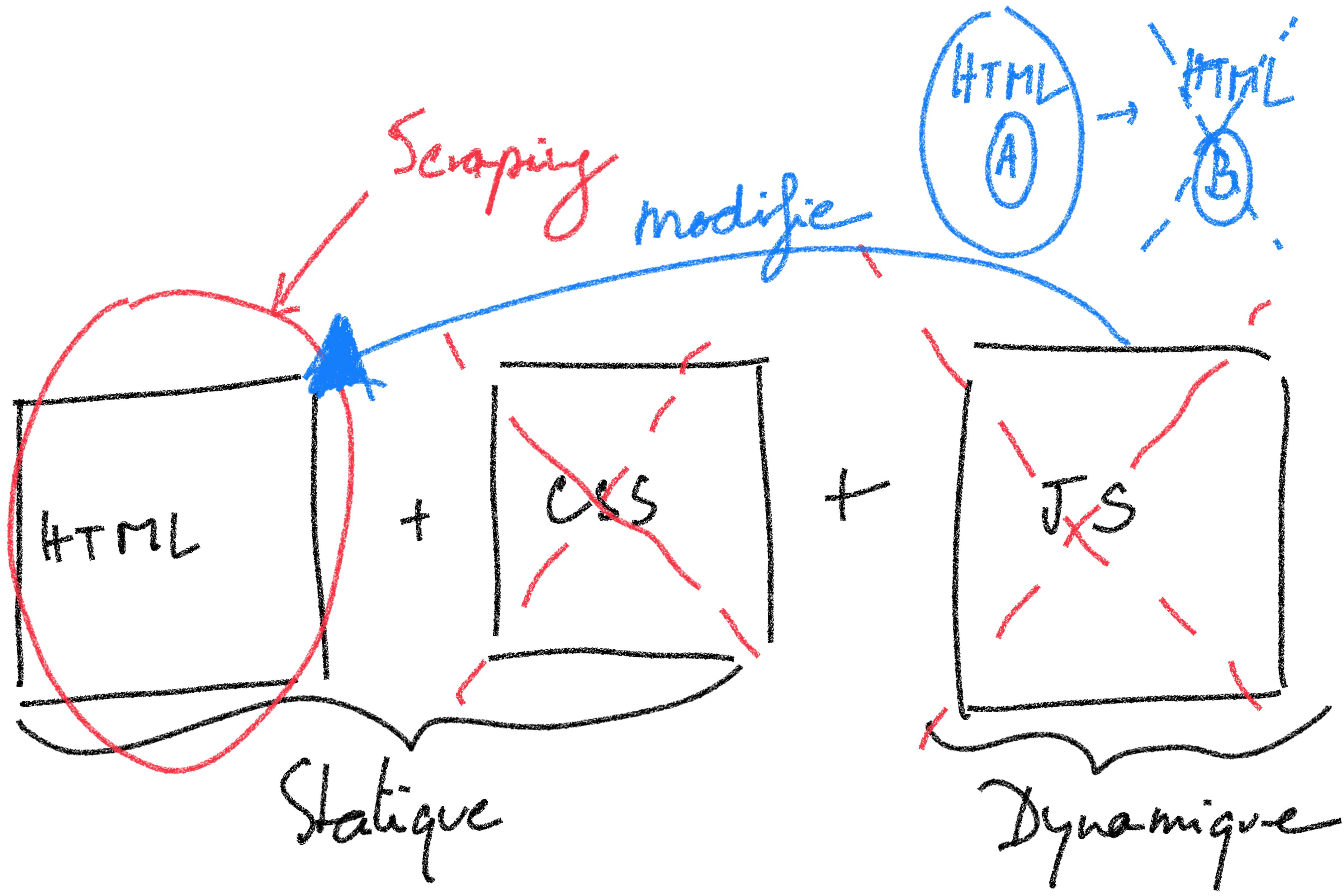
robots.txt

Ferme d'IP

### ②b) Limites de scalabilité :

- i) changer de manière aléatoire les ids / clés des balises (leboncoin.fr)

ii) mettre l'essentiel du contenu de la page dans le JavaScript (et non dans le HTML)



Scraping ++:

Headless Browsers (Selenium)

brows = ChromeBrowser()

brows.go("a/b/c")