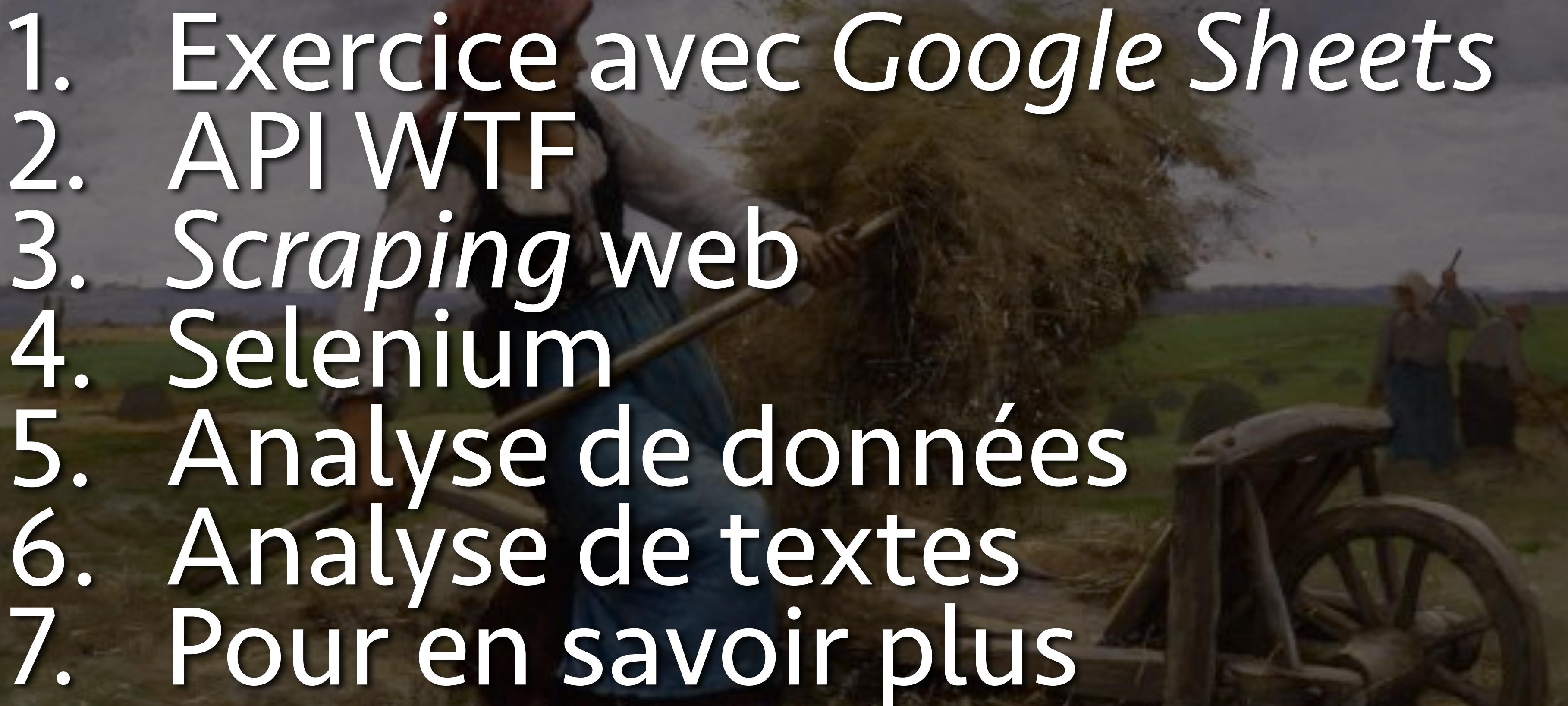


Moissonnage



Moissonnage

- 
1. Exercice avec Google Sheets
 2. API WTF
 3. *Scraping* web
 4. Selenium
 5. Analyse de données
 6. Analyse de textes
 7. Pour en savoir plus

Google Sheets

Fonctions uniques

N'existent pas dans *OO, LO, Excel, Numbers*



=IMPORTHTML



Jean-Hugues Roy



Fil d'actualités

...



Messenger

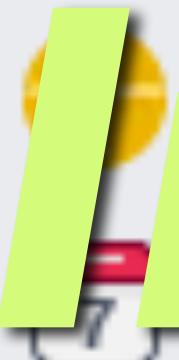


Watch



Marketplace

Raccourcis

 Atelier de journalisme 9 Journées UQTR
 Journées UQTR Mon Fan Club

▼ Voir plus...

Explorer



Interface pour humains

«Qui a décidé que la Terre est ronde?» - Nathalie Lemieux

Dans une publication récente, sur Facebook, la conseillère Nathalie...

```
1 "data": [
2 {
3   "message": "La nomination d'une proche de la Coalition avenir Québec (CAQ) à la tête de la délégation générale du Québec à New York n'est pas partisane, a raconté l'éditorialiste de l'hebdomadaire le Soleil, dans une entrevue à l'agence de presse AP. La nomination d'un membre de la CAQ à la tête de la délégation générale du Québec à New York n'est pas partisane, a raconté l'éditorialiste de l'hebdomadaire le Soleil, dans une entrevue à l'agence de presse AP.", "caption": "lesoleil.com", "status_type": "shared_story", "id": "152738728079652_2276881095665394"
4 },
5 {
6   "message": "Un homme a été enlevé par cinq individus dans le rang Saint-Marc, à Saint-Ambroise, et a été battu sévèrement pour une dette de drogue.", "caption": "lequotidien.com", "status_type": "shared_story", "id": "152738728079652_2276866562333514"
7 },
8 {
9   "message": "MONTRÉAL — Félix Auger-Aliassime et Leylah Annie Fernandez ont lancé leur saison 2019 de brillante façon. Les deux jeunes joueurs de tennis canadiens ont terminé la semaine dernière au deuxième rang mondial de leur discipline.", "caption": "lesoleil.com", "status_type": "shared_story", "id": "152738728079652_2276830112337159"
10 },
11 {
12   "message": "OTTAWA — Les frais de garderie ont baissé — ou ont peu augmenté — dans certaines villes canadiennes, ce qui pourrait être un signe que l'argent dépensé pour l'éducation des enfants est de plus en plus cher.", "caption": "lesoleil.com", "status_type": "shared_story", "shares": {"count": 1}, "id": "152738728079652_2276800845673419"
13 },
14 {
15   "comments": [
16     {
17       "date": [
18         {
19           "message": "Elle est tellement bonne, c'est dommage.", "id": "2276771809009656_2276790839007753"
20         },
21         {
22           "message": "Mme Bédard va nous manquer. Merci pour les beaux reportages.", "id": "2276771809009656_2276782459008591"
23         }
24       ],
25     }
26   ]
27 }
```

Interface pour ordinateurs

API

- Facebook
- Twitter
- Google
 - Drive
 - Maps (\$)
 - Search (\$)
 - Youtube...
- Instagram...
- ~~WhatsApp~~
- Twitch...
- Spotify...
- Uber...
- ~~AirBnb~~
- ...



Application under review.

Thanks! We've received your request for API access and are in the process of reviewing it.

Be sure to watch the [email address associated with this Twitter account](#) at the time of application, as we may request more information to facilitate the review process in the coming days (be sure to check your spam folder as well).

To help us understand how you use your existing apps, please [edit each of your apps](#) and add a description of your app's use case where it says "Tell us how this app will be used".

We know that this application process delays getting started with Twitter's APIs. This information helps us protect our platform and serve the health of the public conversation on Twitter. It also informs product investments and helps us better support our developer community. For more information about our policies please see our [Terms of Service](#) and our [Developer Terms](#).

You'll receive an email when the review is complete. In the meantime, check out our [documentation](#), explore our [tutorials](#), or check out our [community forums](#).

API Twitter

D'abord, se créer une «app»

Intégrer les permissions dans
un script.

api-twitter.py

Utiliser plusieurs mots/expressions

Répéter recherches avec «cron»

Enr. résultats dans base de données

Mais il y a des limites...



Resource Information

Response formats	JSON
Requires authentication?	Yes
Rate limited?	Yes
Requests / 15-min window (user auth)	180
Requests / 15-min window (app auth)	450

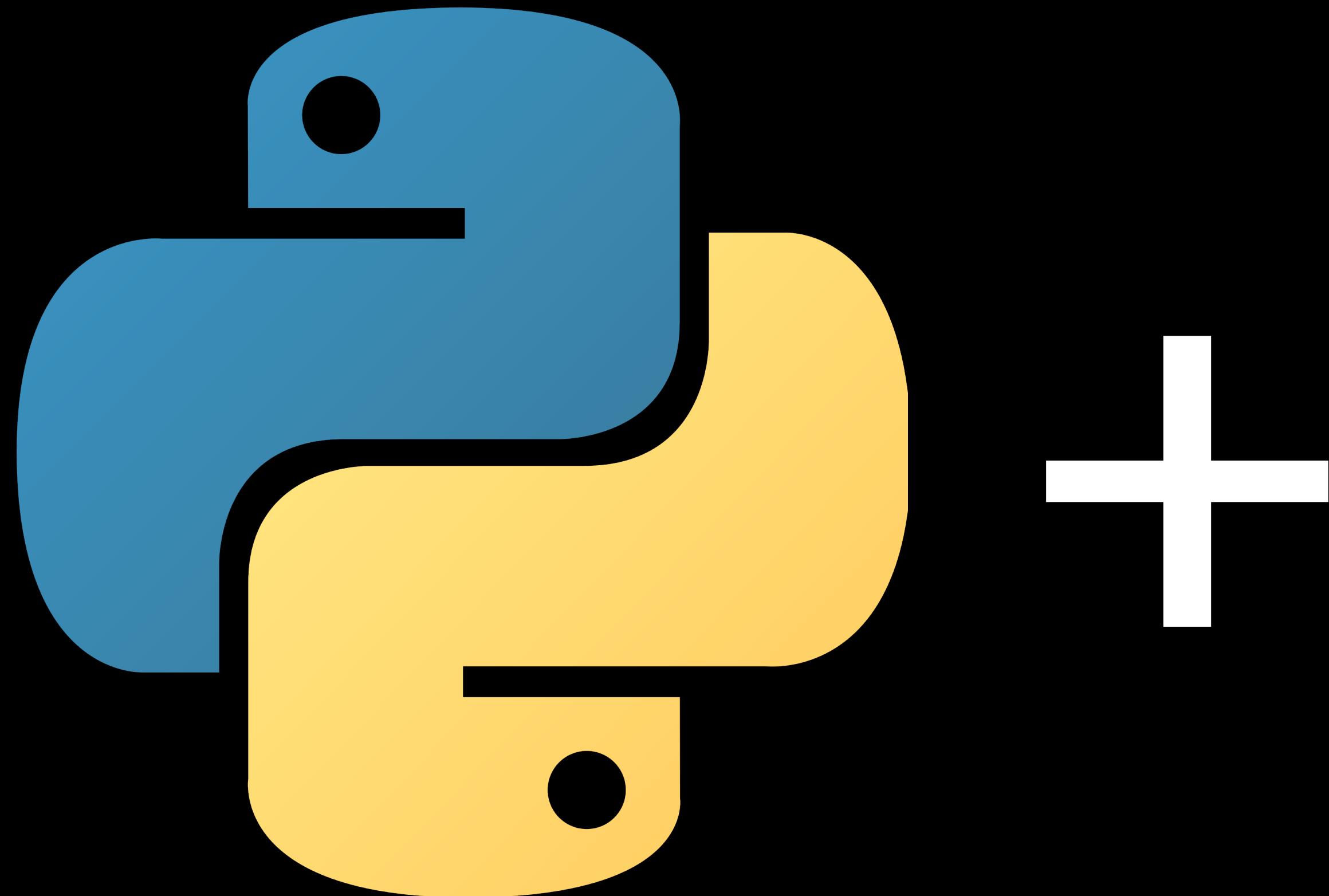


tter.py

ons

ées

Scraping



Python



BeautifulSoup

Scraping

Exemple 1



Objectif :

Ramasser le texte de toutes les lois du Québec en français et en anglais

1^{re} étape :

Recueillir les URLs des lois

lois01.py

Scraping

Exemple 1



2e étape :
Télécharger les 1042 lois (521 dans
chaque langue) `lois02.py`

Fichiers HTML (pas PDF 💩)

Scraping

Exemple 1



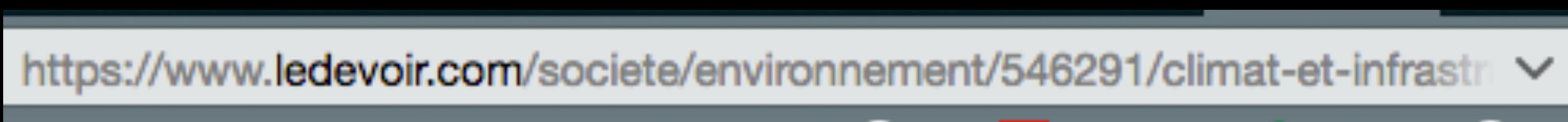
2e étape :
Télécharger les 1042 lois (521 dans
chaque langue) `lois02.py`

Fichiers HTML (pas PDF 💩)

tutoriel

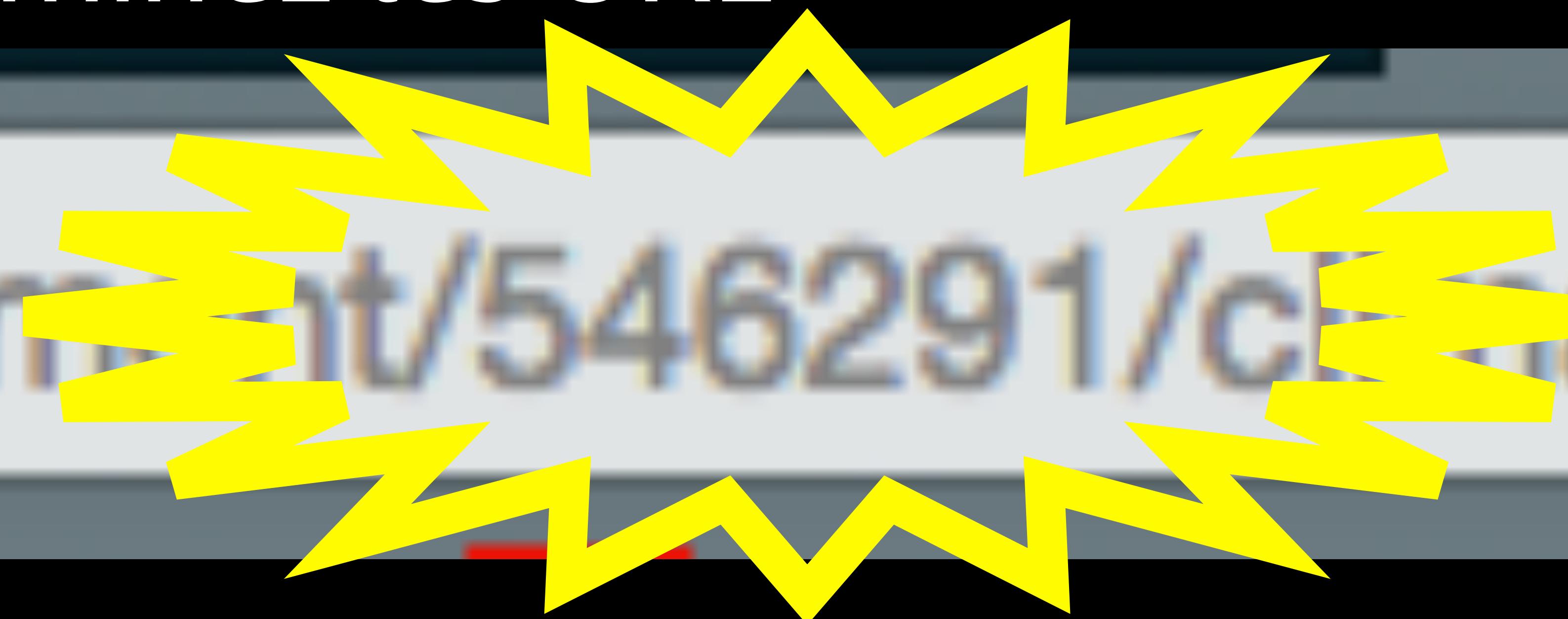
Scraping

Conseils :
Examinez les URL



Scraping

Conseils :
Examinez les URL



Scraping

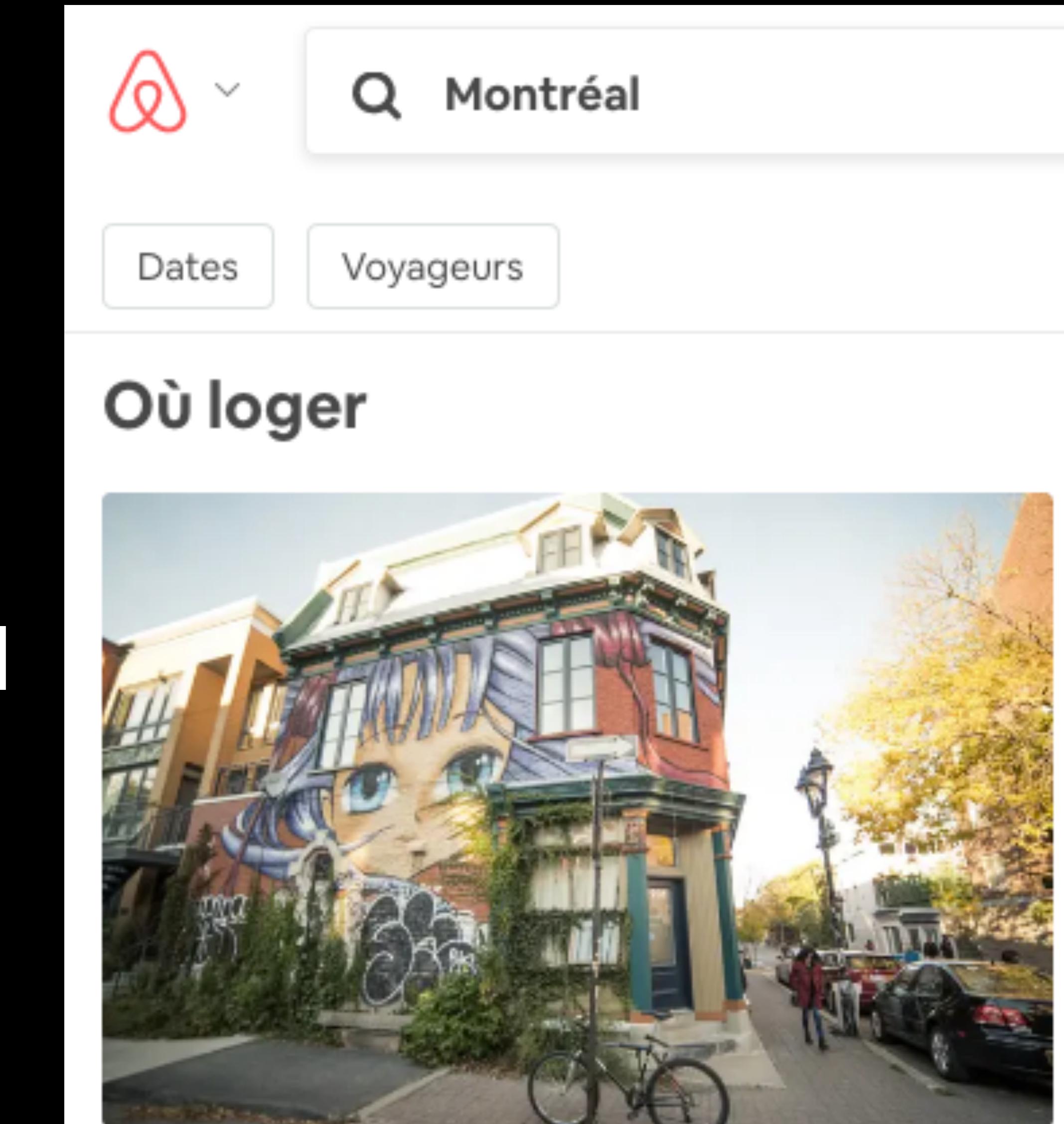
Conseils :

Dans le code HTML,
examinez les balises <meta>

The New York Times

Scraping

Conseils :
Dans le code HTML,
examinez le contenu
de certains scripts



LOFT ENTIER · 1 LIT

LOFT artistic in the heart plateau

\$50 CAD par nuit · Annulation gratuite

★★★★★ 151 · Superhost

Scraping

Stay tuned

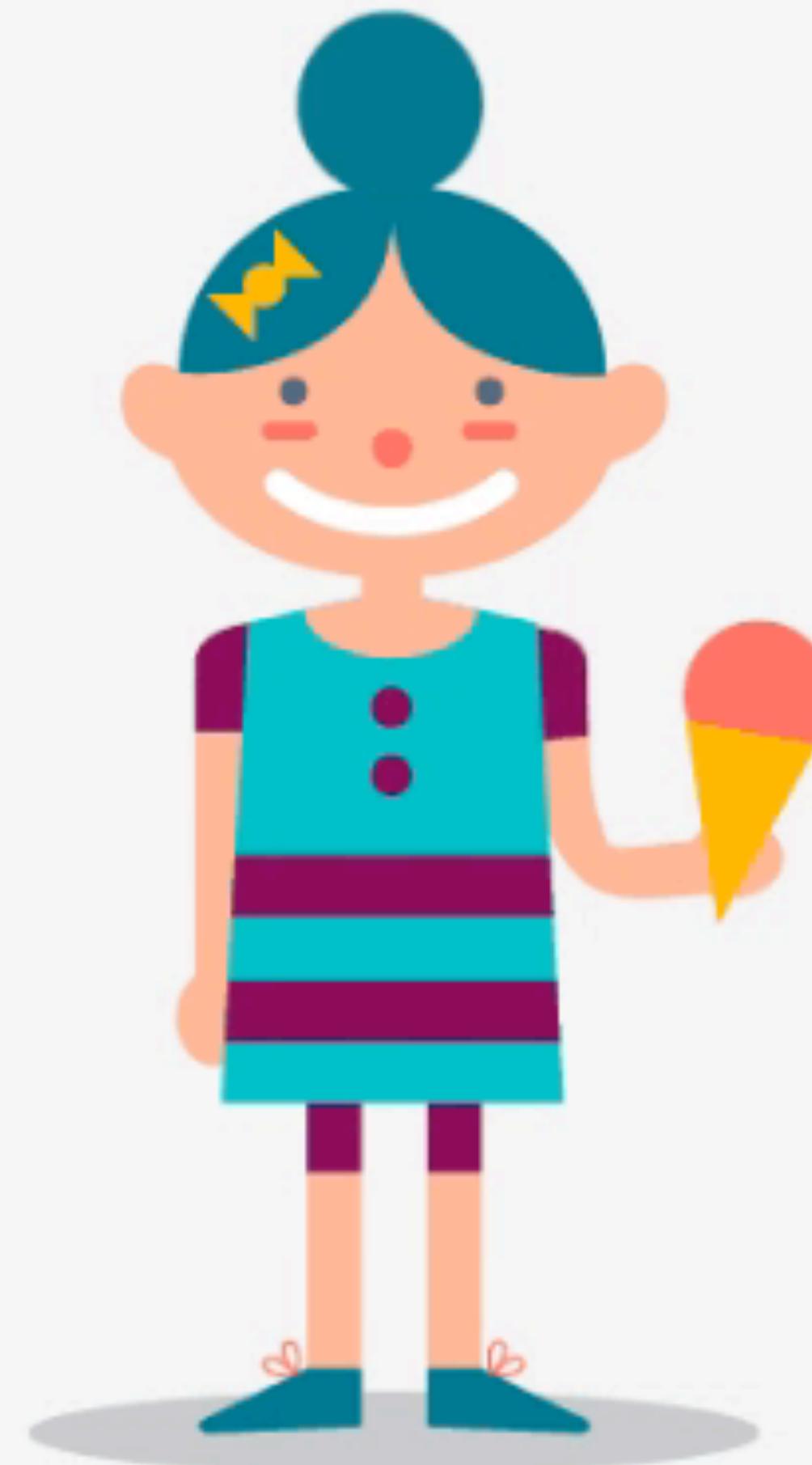
Error code: 503

Airbnb is temporarily unavailable, but we're working hard to fix the problem.
We'll be up and running soon! Keep an eye on our [Twitter account](#) for
updates.

If you need help with an ongoing reservation or for urgent issues, tweet us
[@AirbnbHelp](#) or [call us](#).

Please note, during site downtime, our response times may be longer than
usual.

Thanks for your patience!



Scraping

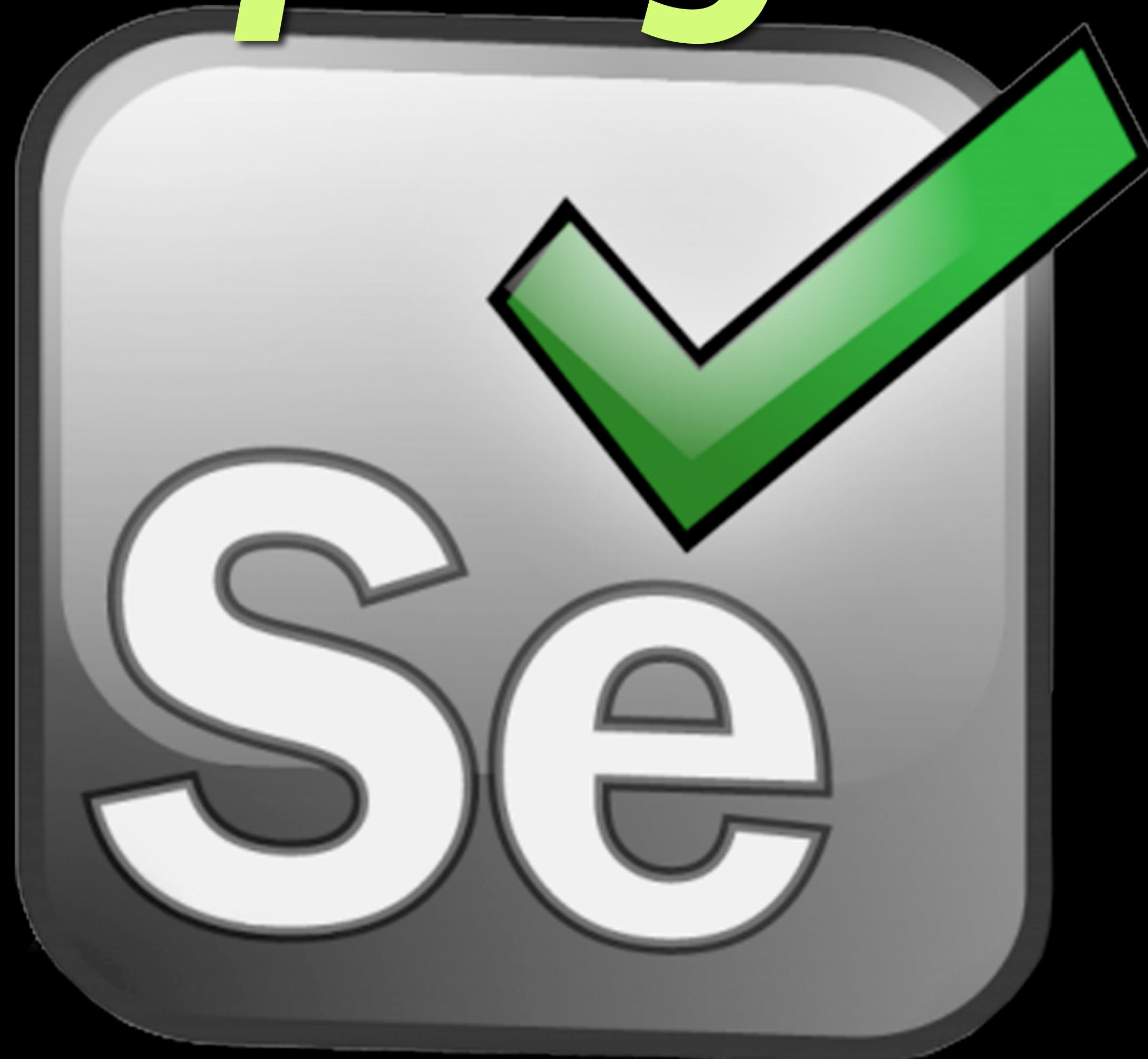
Difficultés :

Et on fait quoi quand ce qu'on cherche
du contenu issu du « deep web »?



COLLÈGE DES MÉDECINS
DU QUÉBEC

Scraping avec Selenium



md.py
bachir.py

Analyse de données

Tableurs pour **chercheuse.eur.s**



Excel



Calc



Numbers



Feuilles de calcul
Google

Analyse de données

The screenshot shows a LibreOffice Calc spreadsheet titled "EDM4434-Tableurs2.ods". The "Éditeur de rapports" (Report Editor) dialog box is open, specifically the "Lignes" (Lines) section. A tooltip box with a dark gray background and white text is overlaid on the left side of the dialog, pointing towards the "Ajouter un champ" (Add field) button. The tooltip contains the following text:

Dans la section **Lignes**, vous cliquez d'abord sur «**Ajouter un champ**» et vous sélectionnez la variable selon laquelle vous souhaitez effectuer votre regroupement.

The "Éditeur de rapports" dialog lists the following fields:

- No du puits
- Nom du puits
- opérateur
- Année
- Municipalité
- Latitude
- Longitude
- profondeur
- unité

The "Municipalité" field is currently selected, highlighted with a light gray background.

Analyse de données



Ouvrir avec LibreOffice



scrapping-nettoye.csv

687,8 Mo

Dernière modification le 2017-08-28 22:14:43

Analyse de données



Analyse de données pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



jupyter



jupyter
notebook

Analyse de textes

Traitement du langage naturel
nltk

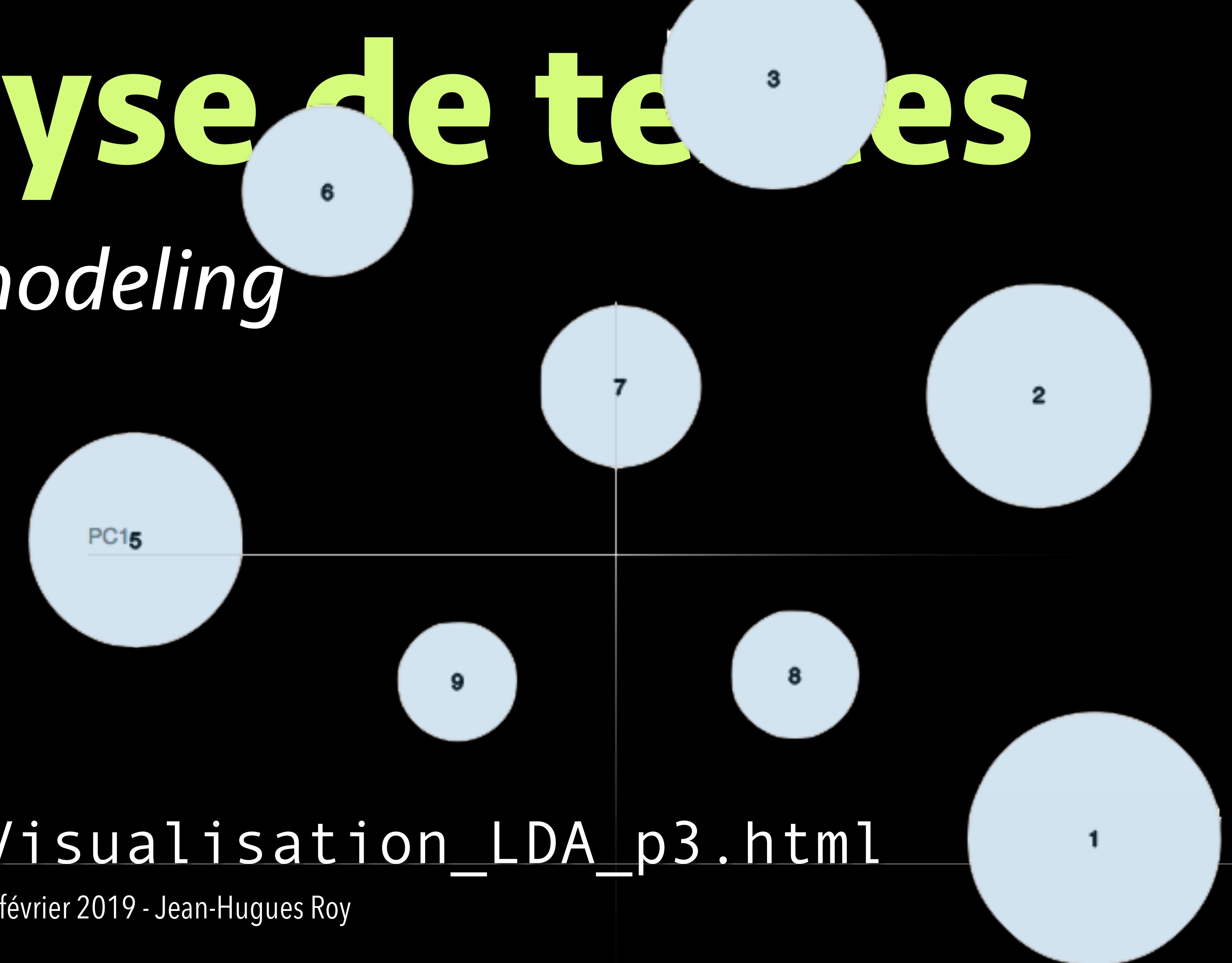
Trois opérations :

- *Tokenization*
- Traitement des mots-vides
- Lemmatisation



Analyse de topics

Topic modeling



Visualisation_LDA_p3.html

Ça vous tente?



ANACONDA

anaconda.com

Ça vous tente?



« Écode l'été » :

- Programmation
- Analyse de données massives
- Apprentissage automatique

2019? Sinon 2020.

Merci!

bit.ly/
abcmojhroy