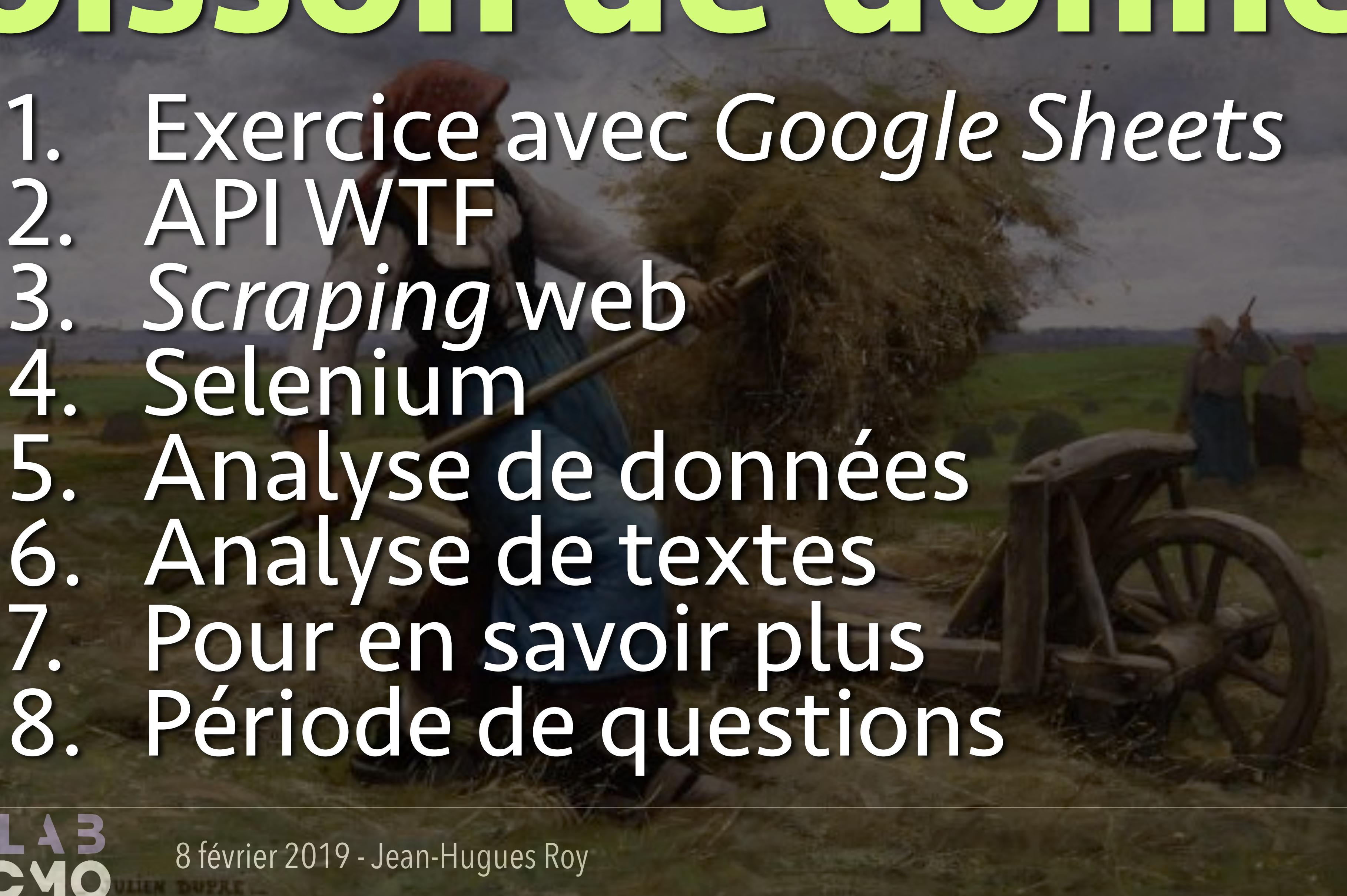


# Moisson de données



# Moisson de données

- 
- A painting of a man harvesting wheat with a scythe in a field.
- 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
  - 8. Période de questions

# Google Sheets

Fonctions uniques

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



=IMPORTHTML



Jean-Hugues Roy



Fil d'actualités

...



Messenger

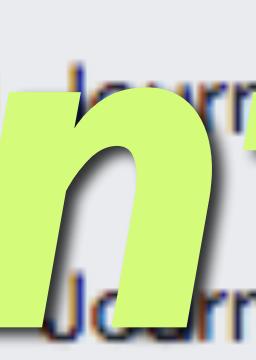


Watch



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 a
4   "caption": "lesoleil.com",
5   "status_type": "shared_story",
6   "id": "152738728079652_2276881095665394"
7 },
8 {
9   "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.",
10  "caption": "lequotidien.com",
11  "status_type": "shared_story",
12  "id": "152738728079652_2276866562333514"
13 },
14 },
15 {
16   "message": "MONTRÉAL – Félix Auger-Aliassime et Leylah Annie Fernandez ont lancé leur saison 2019 de brillante façon. Les deux jeunes joueurs de tennis cana
17   "caption": "lesoleil.com",
18   "status_type": "shared_story",
19   "id": "152738728079652_2276830112337159"
20 },
21 {
22   "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
23   "caption": "lesoleil.com",
24   "status_type": "shared_story",
25   "shares": {
26     "count": 1
27   },
28   "id": "152738728079652_2276800845673419"
29 },
30 {
31   "comments": [
32     "date": [
33       {
34         "message": "elle est tellement bonne, c'est dommage.",
35         "id": "2276771809009656_2276790839007753"
36       },
37       {
38         "message": "Mme Bédard va nous manquer. Merci pour les beaux reportages.",
39         "id": "2276771809009656_2276782459008591"
40       },
41       {
42         "message": "elle est tellement bonne, c'est dommage.",
43         "id": "2276771809009656_2276790839007753"
44       },
45       {
46         "message": "elle est tellement bonne, c'est dommage.",
47         "id": "2276771809009656_2276790839007753"
48       }
49     ]
50   ]
51 }
```

# Interface pour ordinateurs

# Interface pour ordinateurs

# Interface pour ordinateurs

```
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 ra
4     "caption": "lesoleil.com",
5     "status_type": "shared_story",
6     "id": "152738728079652_2276866562333514"
7   },
8   {
9     "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.",
10    "caption": "lequotidien.com",
11    "status_type": "shared_story",
12    "id": "152738728079652_2276866562333514"
13  },
14  {
15    "message": "MONTRÉAL – Félix Auger-Aliassime et Leylah Annie Fernandez ont lancé leur saison 2019 de brillante façon. Les deux jeunes joueurs de tennis cana
16    "caption": "lesoleil.com",
17    "status_type": "shared_story",
18    "id": "152738728079652_2276830112337159"
19  },
20  {
21    "message": "OTTAWA – Les frais de garderie ont baissé ou ont pu augmenté dans certaines villes canadiennes, ce qui pourrait être un signe que l'argent
22    "caption": "lesoleil.com",
23    "status_type": "shared_story",
24    "shares": {
25      "count": 1
26    },
27    "id": "152738728079652_2276800845673419"
28  },
29  {
30    "comments": {
31      "data": [
32        {
33          "message": "C'est malheureux.... elle excellait dans tous ses reportages, triste nouvelle.",
34          "id": "2276771809009656_2276799199006917"
35        },
36        {
37          "message": "Elle est tellement bonne, c'est dommage.",
38          "id": "2276771809009656_2276790839007753"
39        },
40        {
41          "message": "Mme Bédard va nous manquer. Merci pour les beaux reportages.",
42          "id": "2276771809009656_2276782459008591"
43        },
44        {
45        }
46      ]
47    }
48  }
```

REST API  
.json

# API

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

# API Twitter

D'abord, se créer une «app»





## 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

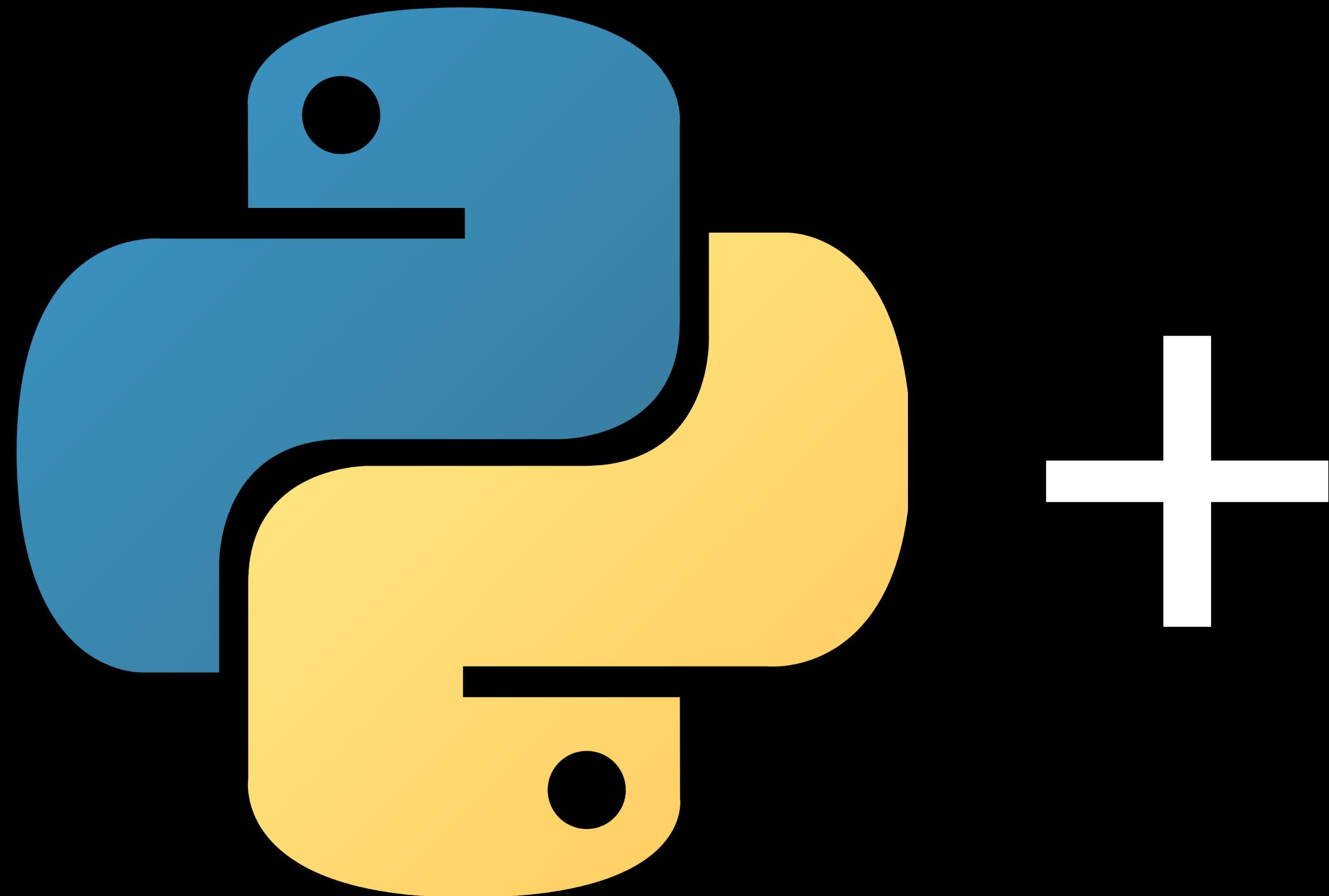


tter.py

ons

ées

# Scraping



Python



BeautifulSoup

# Scraping

Exemple 1



**Objectif :**

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

**1<sup>re</sup> étape :**

Recueillir les URLs des lois

lois01.py

# Scraping

Exemple 1



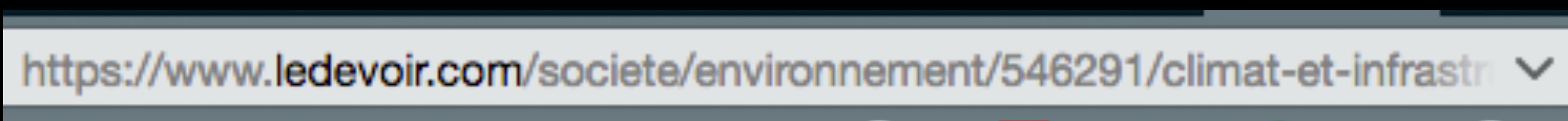
2e étape :  
Télécharger les 1042 lois (521 dans  
chaque langue)

Fichiers HTML (pas PDF 💩) `lois02.py`

tutoriel

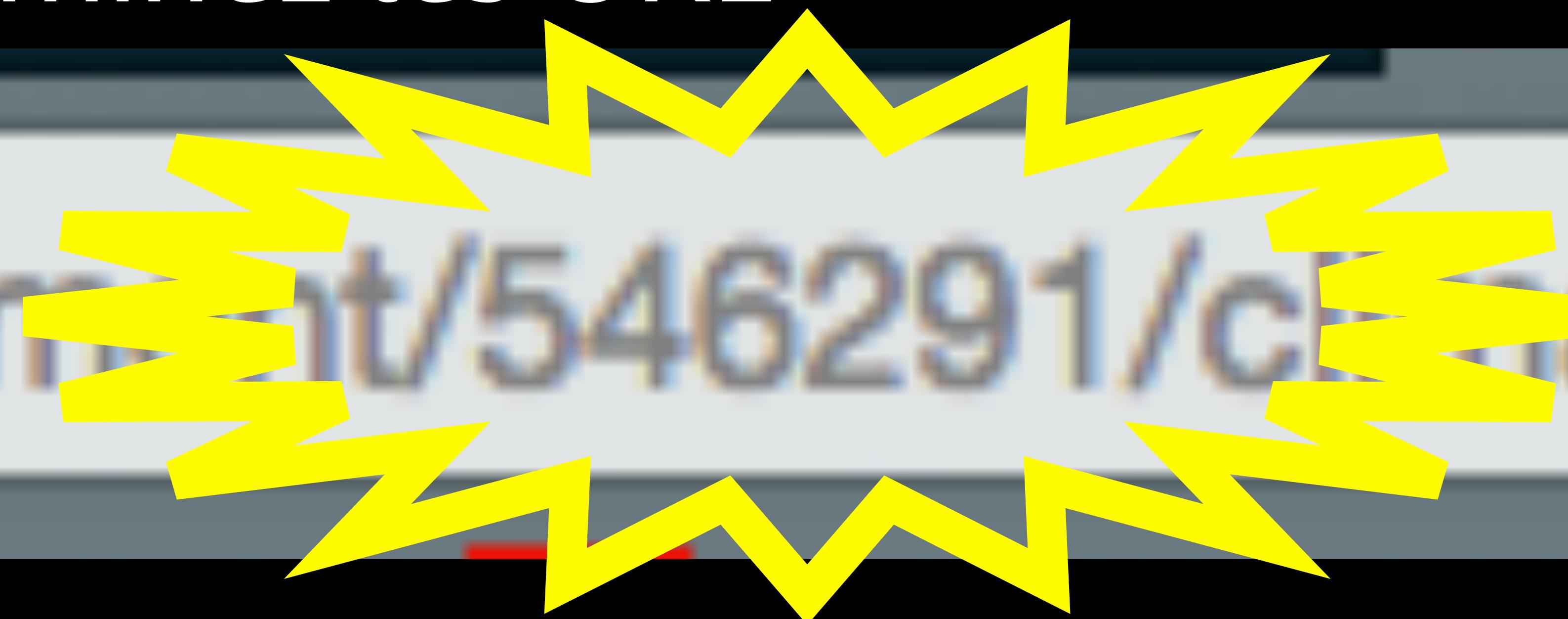
# Scraping

Conseils :  
Examinez les URL



# Scraping

Conseils :  
Examinez les URL



# Scraping

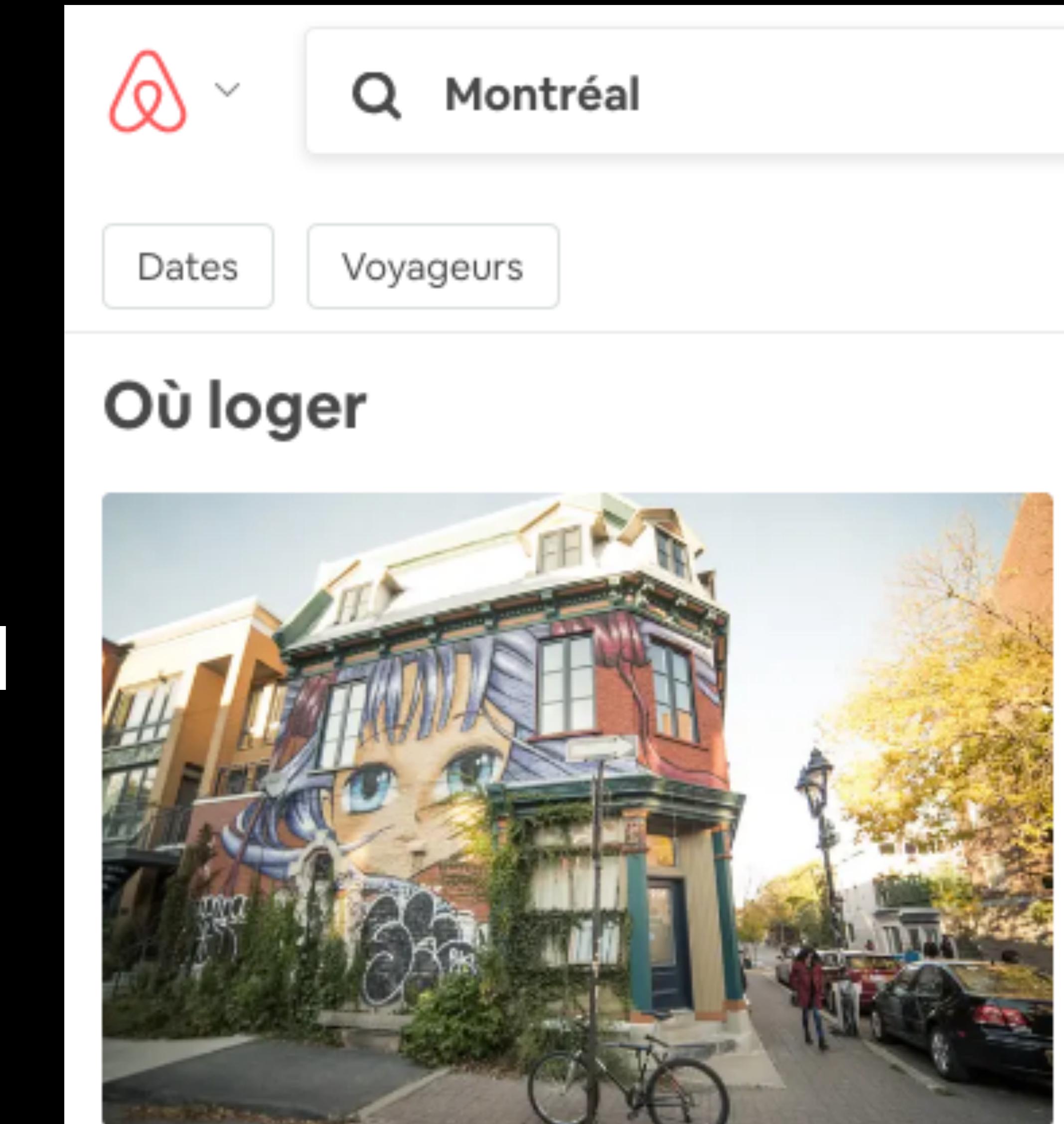
Conseils :

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

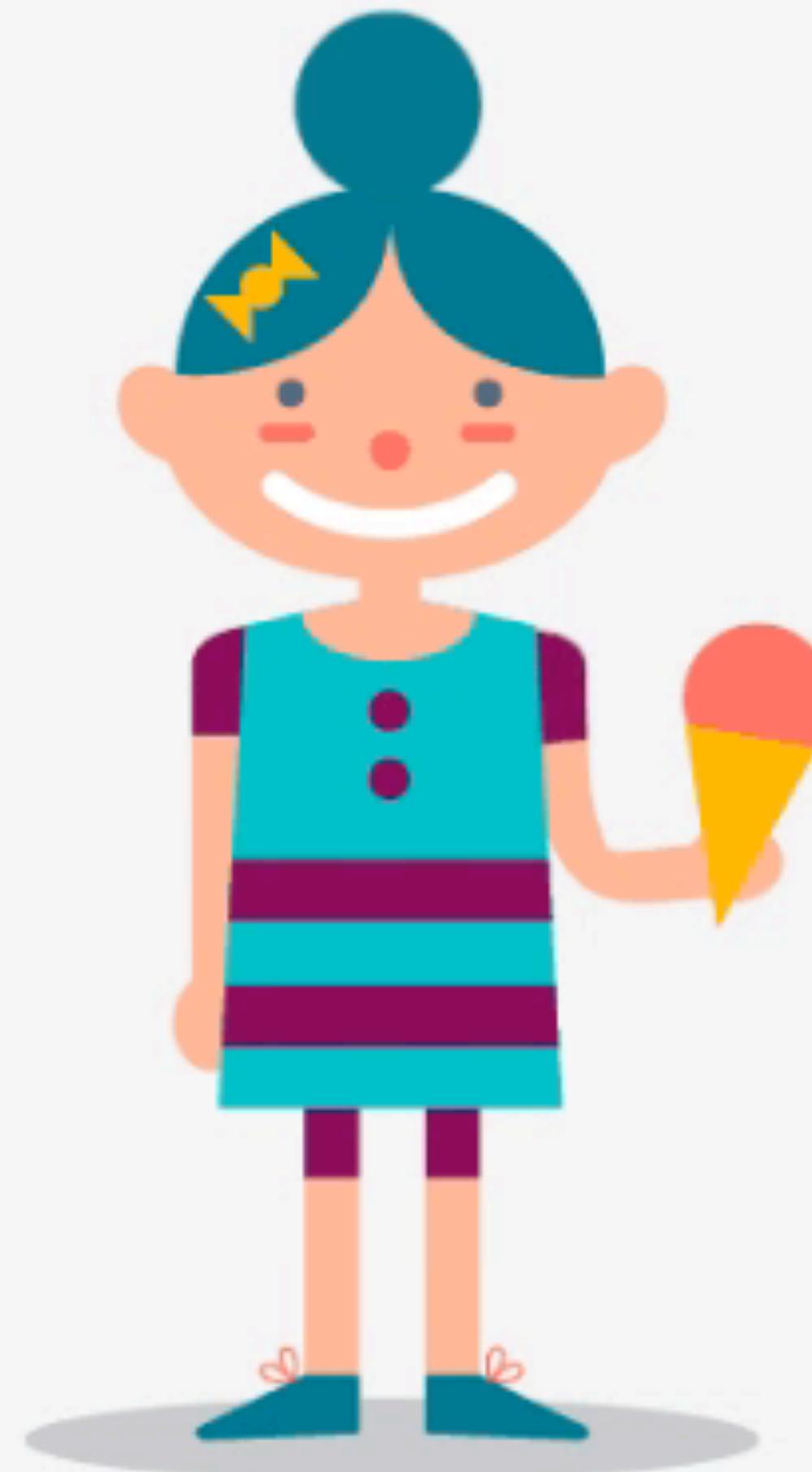
Error code: 503

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

# wc dans Terminal

# Analyse de données





facebook



cole

Sélectionner une base de données

Structure Contenu Relations Déclencheurs Info table Requête

Historique Utilisateurs Connexions

## TABLES

- medias
- medias\_europe
- posts
- posts\_europe
- RCpages
- RCposts

```

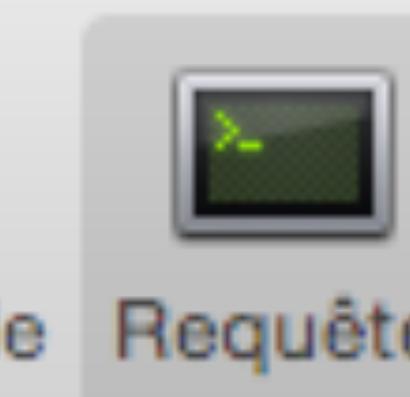
2 sum(partages) as somme_partages,
3 (sum(partages)/count(titre)) as moyenne_partages,
4 sum(reactions) as somme_reactions,
5 (sum(reactions)/count(titre)) as moyenne_reactions,
6 sum(commentaires) as somme_commentaires,
7 sum(likes_commentaires) as somme_likes_commentaires,
8 sum(commentaires_commentaires) as somme_commentaires_commentaires,
9 ((sum(commentaires)+sum(likes_commentaires)+sum(commentaires_commentaires))/count(titre)) as moyenne_commentaires,
10 (sum(partages)+sum(reactions)+(sum(commentaires)+sum(likes_commentaires)+sum(commentaires_commentaires))) as engagement_total,
11 ((sum(partages)+sum(reactions)+(sum(commentaires)+sum(likes_commentaires)+sum(commentaires_commentaires)))/count(titre)) as engagement_moyen,
12 left(dateqc,7) as mois
13 from posts
14 where left(dateqc,4) = "2018"
15 group by titre, mois
16 order by engagement_moyen asc

```

	Favoris	Historique	publications	somme_partages	moyenne_partages	somme_reactions	moyenne_reactions	somme_commentaires	somme_likes_commentaires	Exécuter la requête
		titre								
		Le Journal de Québec	130	6254	48.1077	10661	82.0077	4147	2018-03-04	
		URBANIA	217	11883	54.7604	28053	129.2765	3819	2018-03-04	
		Protégez-Vous	46	6582	143.0870	2648	57.5652	667	2018-03-04	
		CTV Montreal	128	12681	99.0703	9831	76.8047	2147	2018-03-04	
		MétéoMédia	321	23533	73.3115	35819	111.5857	6389	2018-03-04	
		Maurais Live – Radio X	285	17380	60.9825	26432	92.7439	9430	2018-03-04	
		Télé-Québec	77	5850	75.9740	8902	115.6104	1369	2018-03-04	
		Magazine Échos Vedettes	43	565	13.1395	7327	170.3953	2015	2018-03-04	
		Naître et grandir	72	7380	102.5000	6810	94.5833	1596	2018-03-04	
		BLVD 102.1	63	3870	61.4286	4493	71.3175	3503	2018-03-04	
		BLVD 102.1	76	3941	51.8553	4674	61.5000	6297	2018-03-04	
		RDI Économie	59	7388	125.2203	5004	84.8136	701	2018-03-04	
		TVA Gatineau-Ottawa	264	33903	128.4205	14214	53.8409	6392	2018-03-04	
		RDI Économie	91	9155	100.6044	8493	93.3297	1788	2018-03-04	
		TVA Sports	590	8222	13.9356	91105	154.4153	15653	2018-03-04	
		Jeff Fillion	20	1072	53.6000	2241	112.0500	513	2018-03-04	
		Le Journal de Québec	537	25918	48.2644	46783	87.1192	22265	2018-03-04	
		CTV Montreal	74	5588	75.5135	7961	107.5811	1795	2018-03-04	

## INFORMATIONS SUR LA TABLE

- créée : 2018-03-04
- moteur : InnoDB
- lignes : ~546 211
- taille : 366,0 MiB
- encodage : utf8mb4



Structure Contenu Relations Déclencheurs Info table Requête

```

2 sum(partages) as somme_partages,
3 (sum(partages)/count(titre)) as moyenne_partages,
4 sum(reactions) as somme_reactions,
5 (sum(reactions)/count(titre)) as moyenne_reactions,
6 sum(commentaires) as somme_commentaires,
7 sum(likes_commentaires) as somme_likes_commentaires,
8 sum(commentaires_commentaires) as somme_commentaires_commentaires,
9 ((sum(commentaires)+sum(likes_commentaires)+sum(commentaires_commentaires))/count(titre)) as
10 (sum(partages)+sum(reactions)+(sum(commentaires)+sum(likes_commentaires)+sum(commentaires_c
11 ((sum(partages)+sum(reactions)+(sum(commentaires)+sum(likes_commentaires)+sum(commentaires_
engagement_moyen,
12 left(dateqc,7) as mois
13 from posts
14 where left(dateqc,4) = "2018"
15 group by titre, mois
16 order by engagement_moyen asc

```

	Favoris	Historique	publications	somme_partages	moyenne_partages	somme_reactions	moyenne_react
titre							
Le Journal de Québec			130	6254	48.1077	10661	82
URBANIA			217	11883	54.7604	28053	129

# Analyse de données pandas

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



jupyter



jupyter  
notebook

# Analyse de textes

Traitement du langage naturel  
*nltk*

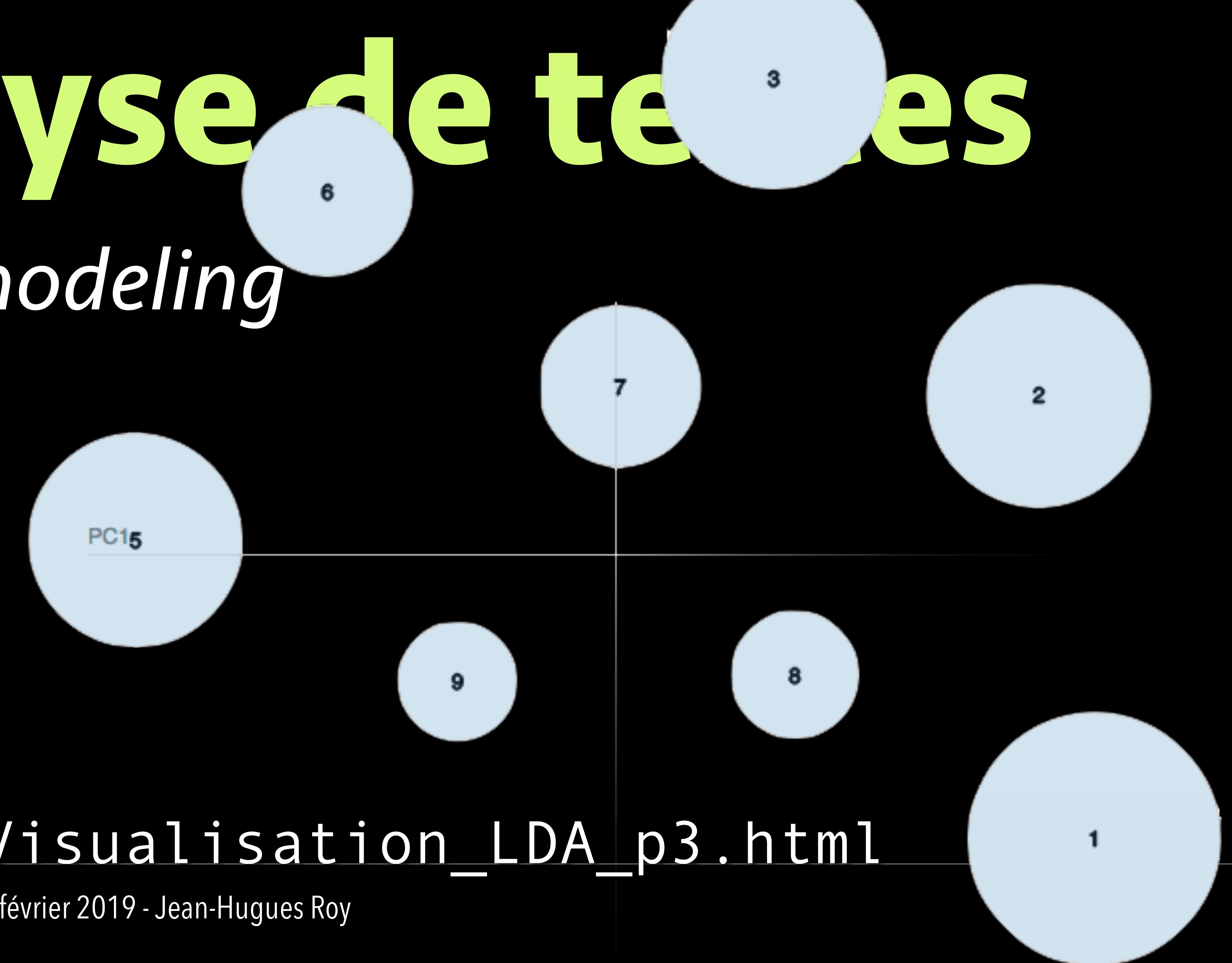
Trois opérations :

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



# Analyse de topics

*Topic modeling*



Visualisation\_LDA\_p3.html

# Ça vous tente?



# ANACONDA

[anaconda.com](https://anaconda.com)

# Ça vous tente?



« Écode l'été » :

- Programmation
- Analyse de données massives
- Apprentissage automatique

2019? Sinon 2020.

# Merci!

bit.ly/  
abcmojhroy