



médialab Tools

De petits outils simples et complexes au service de la DataScience

Benjamin Ooghe-Tabanou
Sciences Po médialab

Sciences Po MÉDIALAB

Collecter et produire des données pour la recherche en SHS Fréjus, 15-18 novembre 2016

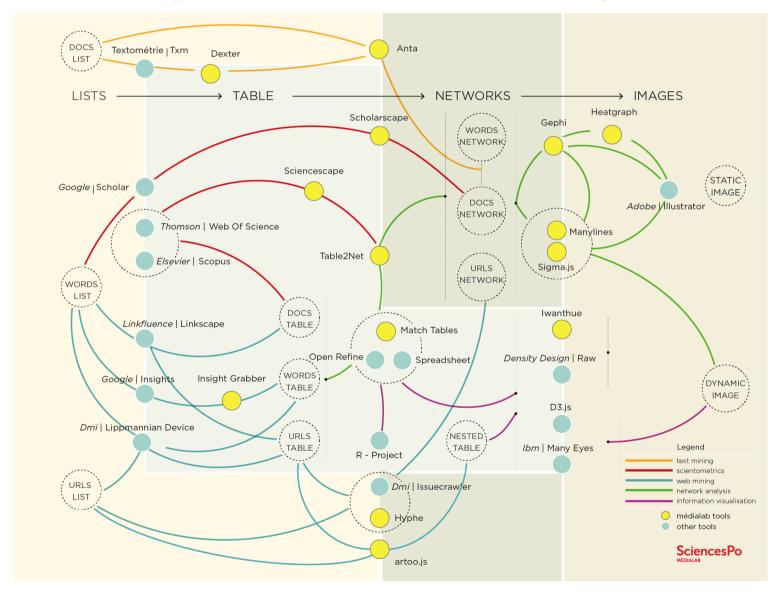






Un écosystème de petits outils pour la science

http://tools.medialab.sciences-po.fr



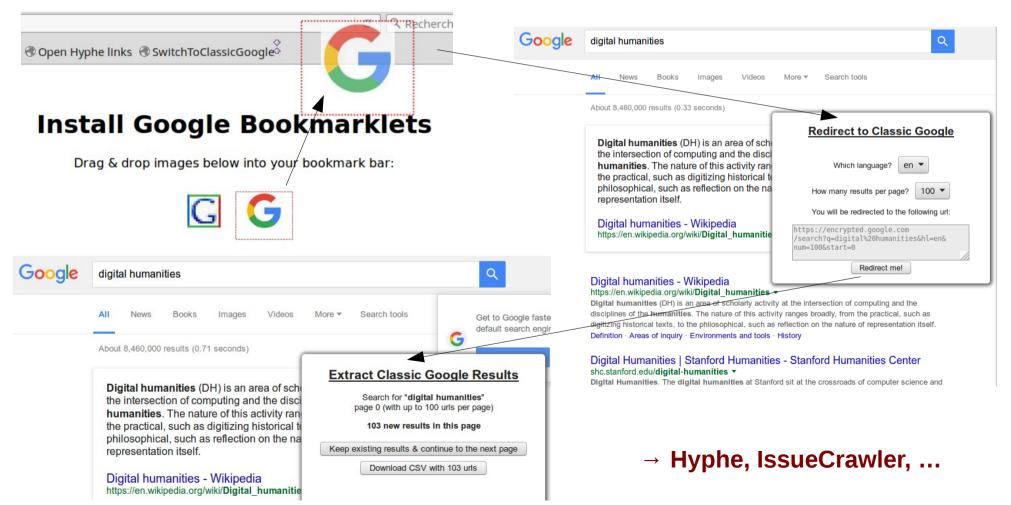




Google bookmarklets : les résultats en CSV

https://medialab.github.io/google-bookmarklets/

• Des boutons dans vos favoris pour récupérer simplement au format tableur les résultats d'une recherche Google



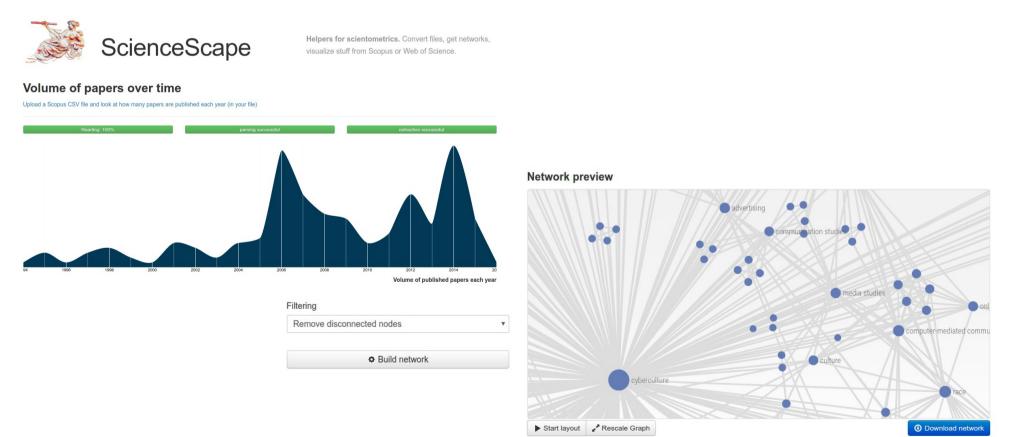




ScienceScape : scientométrie en quelques clics

http://tools.medialab.sciences-po.fr/sciencescape/

- Étudier les auteurs, mots-clés et revues d'un ensemble de publications exportées depuis Scopus ou WebOfScience
 - → exemple : http://jiminy.medialab.sciences-po.fr/data/tools-demo/scopus.csv



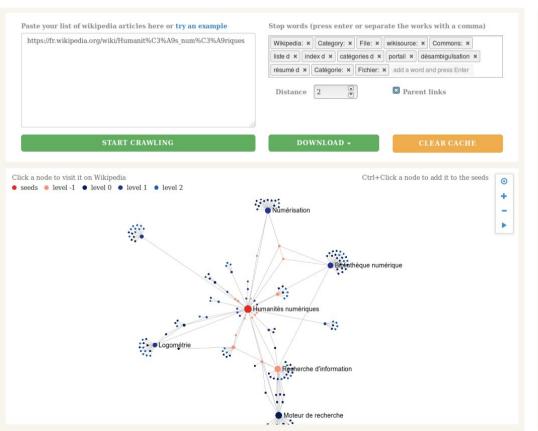


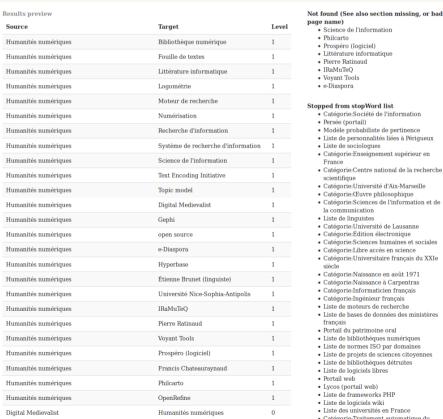


SeeAlsology: exploration sémantique rapide

http://tools.medialab.sciences-po.fr/seealsology/

- Explorer le réseau des liens présents dans les sections « Voir aussi », « Articles connexes » des pages Wikipedia
 - → exemple: https://fr.wikipedia.org/wiki/Humanit%C3%A9s_num%C3%A9riques





 Catégorie:Universitaire français du XXIe Catégorie:Naissance en août 1971 • Catégorie:Naissance à Carpentras Catégorie:Informaticien français Catégorie:Ingénieur français · Liste de moteurs de recherche Liste de bases de données des ministères Liste de bibliothèques numériques · Liste de normes ISO par domaines · Liste de projets de sciences citoyennes Liste de bibliothèques détruites · Liste de logiciels libres · Portail web · Lycos (portail web) • Liste de frameworks PHP

Science de l'information

Catégorie:Société de l'information

· Modèle probabiliste de pertinence

Liste de personnalités liées à Périqueux

· Catégorie:Enseignement supérieur en

Catégorie:Université d'Aix-Marseille

· Catégorie:Œuvre philosophique · Catégorie:Sciences de l'information et de

Catégorie:Université de Lausanne

 Catégorie:Édition électronique Catégorie:Sciences humaines et sociales

Catégorie:Libre accès en science

· Catégorie:Centre national de la recherche

 Prospéro (logiciel) · Littérature informatique

· Pierre Ratinaud IRaMuTeO

Persée (portail)

scientifique

Liste de sociologues

la communication Liste de linguistes

· Liste de logiciels wiki

Liste des universités en France

· Vovant Tools

e.Diaspora

· Philcarto

→ Gephi, Manylines, ...





artoo : extraire des données du web (avancé)

https://medialab.github.io/artoo/

- Un bookmarklet à ajouter dans la barre de favoris du navigateur
- Une librairie JavaScript de fonctions utiles pour le scraping (extraction de données) depuis la console du navigateur (F12)



→ exemple: https://en.wikipedia.org/wiki/List_of_countries_by_carbon_dioxide_emissions

```
> var data = artoo.scrapeTable( ".wikitable", {headers: 'th'} );
undefined
> data.length;
49
> data[0];
Object {Country: " World", CO2 emissions (kt) in 2014[2]: "35,669,000", " %
CO2 Emissions by Country": "100%", Emission per capita (t) in 2014[3]: "5.0"}
> artoo.saveCsv(data, "CO2-world-emissions.csv");
undefined
```

→ CSV-Rinse-Repeat, Table2Net, Khartis, ...





Khartis: cartographier les données d'un CSV

http://www.sciencespo.fr/cartographie/khartis/

- Cartographier des données tabulaires géonommées en quelques clics et exporter une image PNG ou SVG
 - → exemple avec le CSV tiré d'artoo : http://jiminy.medialab.sciences-po.fr/data/tools-demo/CO2-world-emissio







Gazouilloire: extraction de tweets (avancé)

https://github.com/medialab/gazouilloire

- Collecter en direct en continu (et jusqu'à 7 jours en arrière)
 - des tweets par mots-clés, utilisateurs, localisation, langue...
 - les conversations et médias associés
 - des profils d'utilisateurs

```
"user": "Gazou medialab2",
    "oauth_token":
    "oauth secret":
                                          [2016-11-22 15:23:34.056196] DEBUG: Starting search queries with 328 remaining calls for the next 655 seconds
    "host": "localhost",
                                           [2016-11-22 15:23:34.259849] DEBUG: [search] +1 tweets (agriculture%20Paris OR agricultures%20Paris OR agroforesterie%20Paris
    "port": 27017,
                                           2016-11-22 15:23:35.807085] DEBUG: Saved 1 tweets in MongoDB
   "db": "tweets-naturpradi"
                                           [2016-11-22 15:23:37.358533] DEBUG: [search] +1 tweets (espaces%20verts%20Paris OR ferme%20Paris OR fermes%20Paris)
                                           [2016-11-22 15:23:37.810930] DEBUG: Saved 1 tweets in MongoDB
                                           [2016-11-22 15:23:45.049743] DEBUG: [stream] +1 tweet
   "écologique Paris",
                                           2016-11-22 15:23:45.821150] DEBUG: Saved 1 tweets in MongoDB
   "végétation Paris",
                                           2016-11-22 15:24:51.598045] DEBUG: [stream] +1 tweet
   "verger Paris",
    "grenelle environnement Paris",
                                           [2016-11-22 15:24:51.893009] DEBUG: Saved 1 tweets in MongoDB
   "locavore Paris"
                                           2016-11-22 15:24:52.401661] DEBUG: [medias] +1 files
                                           [2016-11-22 15:24:58.073013] DEBUG: Starting search queries with 286 remaining calls for the next 571 seconds
"time_limited_keywords": {
                                           [2016-11-22 15:25:00.383614] DEBUG: [stream] +1 tweet
                                           [2016-11-22 15:25:01.905385] DEBUG: Saved 1 tweets in MongoDB
"geolocalisation": null,
                                           [2016-11-22 15:26:18.060840] DEBUG: Starting search queries with 246 remaining calls for the next 491 seconds
"geolocalisation_type": "admin",
                                          [2016-11-22 15:26:19.922864] DEBUG: [search] +1 tweets (compost%20Paris OR composts%20Paris OR compostage%20Paris)
"resolve_redirected_links": true,
                                          [2016-11-22 15:26:19.989779] DEBUG: Saved 1 tweets in MongoDB
'grab_conversations": true,
"download medias": true,
"medias_directory": "/store/tweets/naturpradi/media/",
"timezone": "Europe/Paris",
"debug": true
```

 \rightarrow exemple d'export CSV: http://jiminy.medialab.sciences-po.fr/data/tools-demo/tweets-data-humanities.csv

→ Catwalk, CSV-Rinse-Repeat, ...

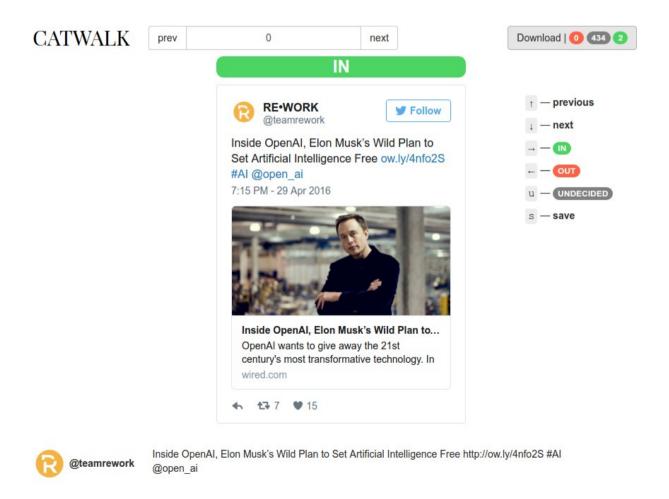




CatWalk : sélection qualitative de tweets

https://medialab.github.io/catwalk/

 Passer en revue rapidement « à la Tinder » tous les tweets d'un CSV pour décider de les inclure / exclure d'un corpus



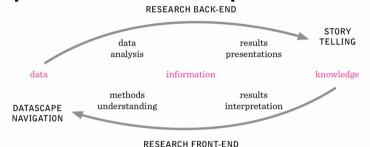




CSV-Rinse-Repeat : exploration de CSV (avancé)

http://tools.medialab.sciences-po.fr/csv-rinse-repeat/

- Itérations successives pour identifier problèmes & questions
- Nettoyer, filtrer, explorer, visualiser, enrichir et exporter en JavaScript le contenu d'un CSV



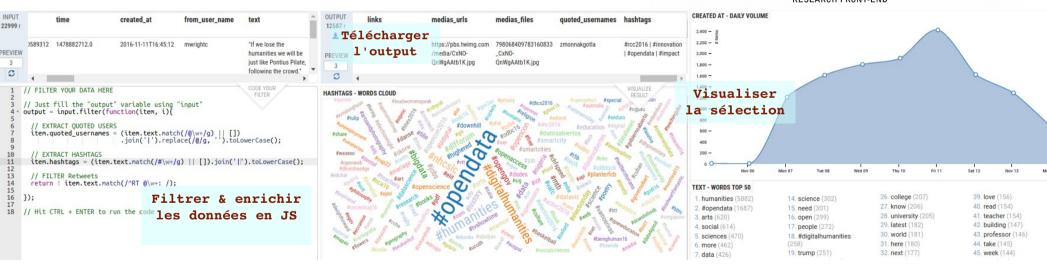


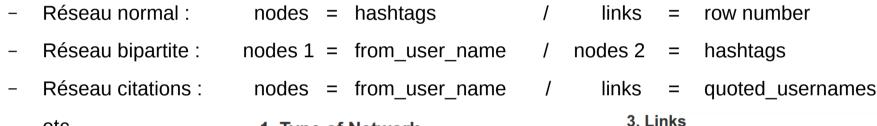




Table2Net : faire un réseau à partir d'un CSV

http://tools.medialab.sciences-po.fr/table2net/

- Générer un graphe de liens entre éléments à partir des données d'un fichier tableur
- Exemples tirés de CSV-Rinse-Repeat: http://jiminy.medialab.sciences-po.fr/data/tools-demo/tweets-data-humanities-rinsed





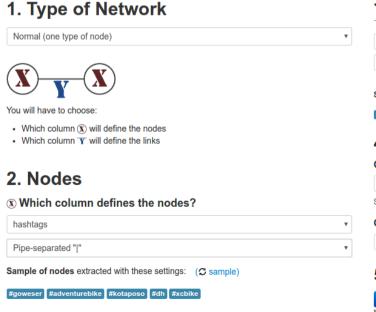


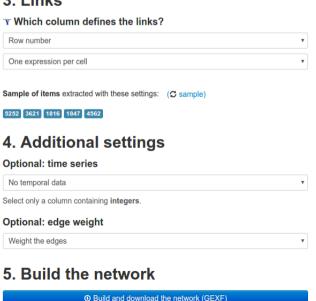
Load your CSV table

It has to be **comma-separated** and the first row must be dedicated to **column names**.

Parsing successful. 44 columns and 5576 rows.

→ Gephi, Manylines, ...





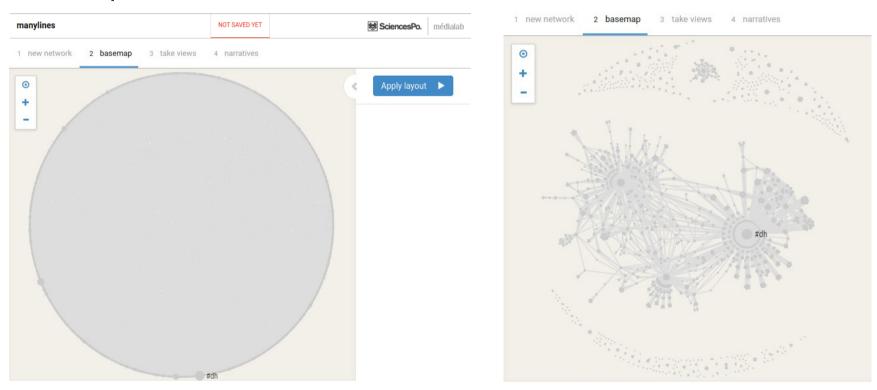




Manylines : publier et documenter un réseau

http://tools.medialab.sciences-po.fr/manylines

- Explorer et publier rapidement en ligne un réseau
- Raconter le réseau sous la forme de slides de présentation
 - → ex exporté de Table2Net : http://jiminy.medialab.sciences-po.fr/data/tools-demo/tweets-data-humanities-hashtags-network.gexf



• Exemple de slides résultants : http://tools.medialab.sciences-po.fr/manylines/embed#/narrative/5aaec64d-96c1-46





Merci de votre attention!

Sciences Po MÉDIALAB

@medialab_ScPo

benjamin.ooghe@sciencespo.fr

