

# Des archives du commerce à des données quantifiables

*une longue chaîne de transformation des données*

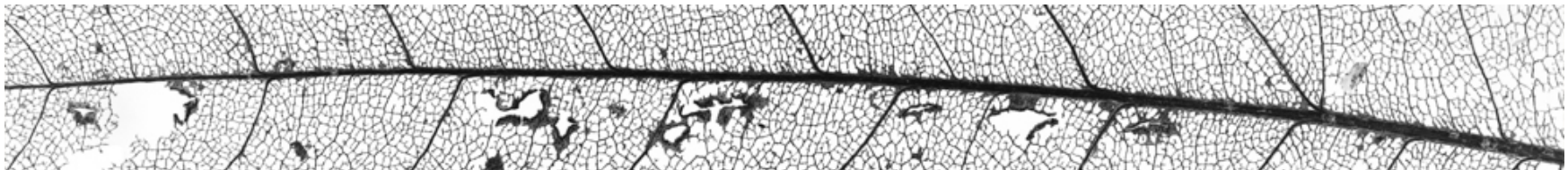
Paul Girard

Sciences Po, médialab

Collecter et produire des données pour la recherche en SHS

Axe 2 : Collecter des données pour les exploiter : comment les préparer en amont ?

Fréjus, le 16 novembre 2016





- Conception de méthodes numériques
- Hybridant les approches qualitatives et quantitatives
- Développant des outils-logiciels
- En Sciences Humaines et Sociales

# Des archives du commerce à des données quantifiables

*une longue chaîne de transformation des données*

- Retour d'expérience des projets:
  - RICardo
  - TOFLIT18
- Des archives aux données ?
- Quelles bases de données ?
- L'exploration visuelle au service des données

Latour, Bruno. 1993. *'Le Topofil de Boa-Vista. La Référence Scientifique: Montage Photophilosophique'*. *Raisons Pratiques* 4: 187–216.

# RICardo · XIXème siècle

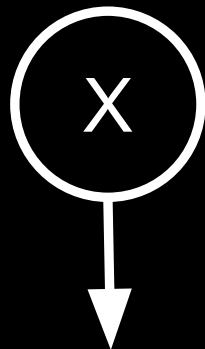
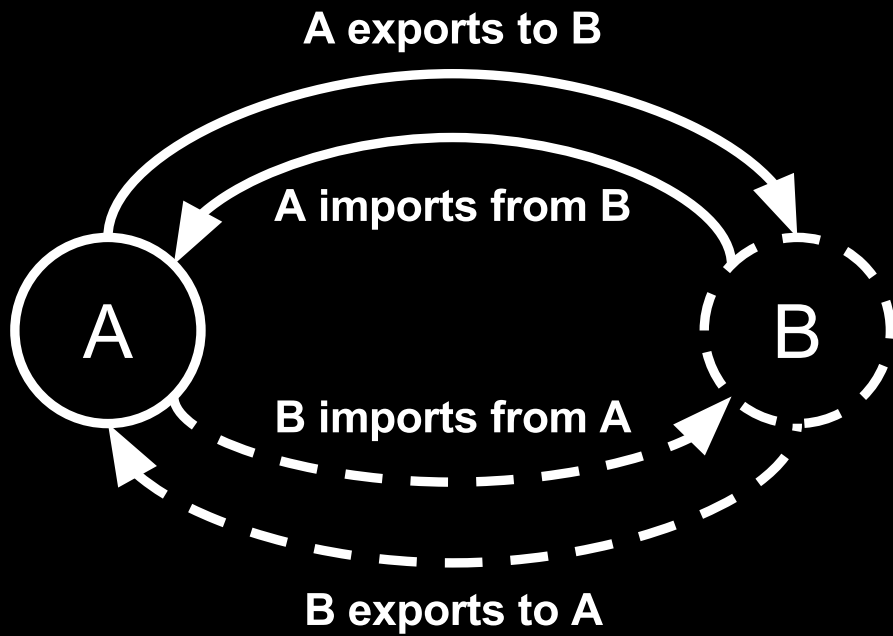
*données annuelles*

1787 · · · · · 1938

# RICardo · commerce bilateral

Flux de commerce entre ~~pays~~ entités

A <··> B



**Political entities**

**trade flow (annual value)  
as reported by X**

# RICardo · commerce total

Flux de commerce entre les entités A et le monde

A <··> M

ANGLETERRE.

La navigation de Liverpool a donné lieu, en 1833, au mouvement ci-après :

LIVERPOOL.

Mouvement commercial en 1833.

TOFLIT18: le XVII<sup>l</sup>ème siècle

données annuelles

1719

1839

PROVENANCES ET DESTINATIONS.	ENTRÉE.		SORTIE.		
	Navires.	Tonneaux.	Navires.	Tonneaux.	
États-Unis.....	562	224,220	408	162,000	
Brésil.....	93	23,550	138	34,500	
Hollande.....	249	27,930	261	29,430	
Belgique.....	14	23,396	15	27,300	
Prusse.....	7	23,500	9	27,700	
Danemarck.....	78	15,600	125	24,600	
Russie.....	68	21,200	84	13,275	
Villes anséatiques.....	60	11,668	59	9,350	
Deux-Siciles.....	83	12,190	38	6,655	
Suède et Norwège.....	42	10,500	24	4,200	
Espagne.....	57	6,840	36	6,480	
Portugal.....	60	7,200	54	6,540	
Indes-Occidentales.....	20	5,248	20	6,100	
France.....	45	6,109	39	4,831	
Plata (États-Unis de la)...	20	3,000	33	7,095	
Mexique.....	24	3,600	33	5,775	
Sardaigne.....	18	2,121	50	6,275	
Colombie.....	18	2,700	12	4,800	
Toscane.....	6	1,200	32	5,760	
Autriche.....	20	3,000	24	3,600	
Chili.....	12	3,000	12	3,000	
Turquie.....	25	2,750	24	3,000	
Autres contrées.....	43	10,020	54	12,075	
Amérique	continentale...	345	120,750	321	112,350
	Antilles.....	230	46,000	201	55,275
Afrique.....	48	14,409	51	17,650	
Indes-Orientales.....	51	12,160	60	11,000	

Pays étrangers

Possessions



# TOFLIT18: commerce de la France

Flux de commerce entre la France et ses partenaires commerciaux

France <··> A,B,C

*rapportés par l'état français*

# TOFLIT18 : les sources

10

Cabas de palme... Espagne

1515 .

Cabarets... Hollande

38 .

Cabillaux { Danemark  
Hollande  
Nord

93108 . 8  
207966 . 18  
156 .

301531 . 6

Cable... Hollande

800 .

Cacav { Espagne  
Flandre  
Hollande  
Isle  
Italie  
Nord  
Portugal  
Savoie

22079 . 16  
1730 .  
500337 . 13  
28687 . 10  
192757 .  
18751 .  
166183 .  
5200 .

830625 . 10

Cachou... Hollande

757 . 10

# TOFLIT18: les produits

## Top 50 des produits

Articles réunis · Indigo · Mercerie · Eau de vie · Librairie ·  
Vinaigre · Cacao · Suif · Beurre · Alun · Liqueurs · Sel · Fromage ·  
Rocou · Confitures · Acier · Fer ; en barres · Thé · Huile d'olive ·  
Porcelaine · Poivre · Farine · Miel · Chandelle · Huile ; d'olive ·  
Fayance · Savon · Amidon · Ris · Verdet · Bière · Cochenille ·  
Légumes · Bijouterie · Garance · Horlogerie · Chocolat · Meubles ·  
Quinquina · Amandes · Crin · Papier ; blanc · Planches ; de sapin ·  
Jambons · Lard · Drogues réunies · Argenterie · Bougie · Gaudron ·  
Cuivre...

# Des sources aux données

Volumes d'archives > images > ? > chercheurs

# Transcription manuelle



266947 rows

Show as: rows records Show: 5 10 25 50 rows

	flow	unit	currency	year	reporting	partner	export_import	special_general	species_b
ier.hit-u.ac.jp/COE/Japanese/discussionpapers/DP97.28/table2.htm	6122.0	1000	pesos	1855	philippines	Total	exp	0	
ier.hit-u.ac.jp/COE/Japanese/discussionpapers/DP97.28/table2.htm	9133.0	1000	pesos	1856	philippines	Total	exp	0	
ier.hit-u.ac.jp/COE/Japanese/discussionpapers/DP97.28/table2.htm	11896.0	1000	pesos	1857	philippines	Total	exp	0	

### Cluster & Edit column "partner"

This feature helps you find groups of different cell values that might be alternative representations of the same thing. For example, the two strings "New York" and "new york" are very likely to refer to the same concept and just have capitalization differences, and "Gödel" and "Godel" probably refer to the same person. [Find out more ...](#)

Method: **key collision** Keying Function: **fingerprint** **390 clusters found**

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
5	576	<ul style="list-style-type: none"> <li>Suède et Norvège (542 rows)</li> <li>Norvège et Suède (17 rows)</li> <li>Suède et Norvège** (12 rows)</li> <li>Suède, Norvège (2 rows)</li> <li>Suède et norvège (2 rows)</li> </ul>	<input type="checkbox"/>	Suède et Norvège
4	26	<ul style="list-style-type: none"> <li>Villes anséatiques (Hambourg) (16 rows)</li> <li>Villes Anséatiques (Hambourg) (7 rows)</li> <li>Villes Anséatiques - Hambourg (2 rows)</li> <li>Villes Anséatiques (Hambourg)** (1 rows)</li> </ul>	<input type="checkbox"/>	Villes anséatiques (Hambourg)
4	106	<ul style="list-style-type: none"> <li>British North Borneo (82 rows)</li> <li>Borneo (British): British North Borneo (16 rows)</li> <li>British north borneo (6 rows)</li> <li>Borneo, British (British North Borneo) (2 rows)</li> </ul>	<input type="checkbox"/>	British North Borneo
4	694	<ul style="list-style-type: none"> <li>Pérou (683 rows)</li> <li>PEROU (4 rows)</li> <li>Perou (4 rows)</li> <li>Pérou** (3 rows)</li> </ul>	<input type="checkbox"/>	Pérou
4	70	<ul style="list-style-type: none"> <li>South Africa, Cape of Good Hope (31 rows)</li> <li>Cape of Good Hope (South Africa) (30 rows)</li> </ul>	<input type="checkbox"/>	South Africa, Cape of Good Hope

# Choices in Cluster

Average Length of Choices

Length Variance of Choices

Rows in Cluster

0 — 9000

0 — 71

0 — 14

Select All Unselect All Merge Selected & Re-Cluster Merge Selected & Close Close

# Diagnostique des données

- Contrôle qualité des données par script
- Analyse quantitative comme aide au nettoyage qualitatif

ier.hit-u.ac.jp/COE/Japanese/discussionpapers/DP97.28/table2.htm	26293.0	1000	pesos	1888	philippines	Total	exp	0
ier.hit-u.ac.jp/COE/Japanese/discussionpapers/DP97.28/table2.htm	34927.0	1000	pesos	1889	philippines	Total	exp	0
ier.hit-u.ac.jp/COE/Japanese/discussionpapers/DP97.28/table2.htm	26214.0	1000	pesos	1890	philippines	Total	exp	0
ier.hit-u.ac.jp/COE/Japanese/discussionpapers/DP97.28/table2.htm	26905.0	1000	pesos	1891	philippines	Total	exp	0
ier.hit-u.ac.jp/COE/Japanese/discussionpapers/DP97.28/table2.htm	27977.0	1000	pesos	1892	philippines	Total	exp	0

# Rapport de test-données de Ricardo

```
# RICNames TEST
0 missing reporting in RICNames
0 missing partners in RICNames
missings written in out_data
# CURRENCY TEST
total number of currencies in flow 12366
check number before/after set currency : 12366/12366
check number before/after set modified_currency : 16382/16382
in currency not in flow 4016
in flow not in currency 0
in flow and in currency 12366
in flow in currency not in rate 1
total known currencies in flow 12365
missign rates exported in out_data
# EXP IMP TEST
missing expimp spe/gen in standards :0
EXP IMP TEST : OK
# FLOWS DUPLICATES TEST
## Spe/gen Dups
```





MOIS del Ere Vulgaire	MOIS Germinal	MOIS del Ere Vulgaire	MOIS Floréal	MOIS del Ere Vulgaire	MOIS Prairial	MOIS del Ere Vulgaire	MOIS Messidor	MOIS del Ere Vulgaire	MOIS Thermidor	MOIS del Ere Vulgaire	MOIS Fructidor
<i>Mars 1794</i>	VII <sup>e</sup> MOIS	<i>Avril 1794</i>	VIII <sup>e</sup> MOIS	<i>Mai 1794</i>	IX <sup>e</sup> MOIS	<i>Juin 1794</i>	X <sup>e</sup> MOIS	<i>Jull. 1794</i>	XI <sup>e</sup> MOIS	<i>Av. 1794</i>	XII <sup>e</sup> MOIS
v. 21 s. 22 D. 23	1 Primidi 2 Duodi 3 Triidi	D. 20 P. 21 L. 22	1 Primidi 2 Duodi 3 Triidi	P. 20 L. 21 M. 22	1 Primidi 2 Duodi 3 Triidi	P. 20 L. 21 M. 22	1 Primidi 2 Duodi 3 Triidi	P. 20 L. 21 M. 22	1 Primidi 2 Duodi 3 Triidi	P. 20 L. 21 M. 22	1 Primidi 2 Duodi 3 Triidi
v. 22 s. 23 D. 24	4 Quaridi 5 Quinidi 6 Sexidi	M. 23 J. 24 V. 25	4 Quaridi 5 Quinidi 6 Sexidi	V. 23 L. 24 M. 25	4 Quaridi 5 Quinidi 6 Sexidi	V. 23 L. 24 M. 25	4 Quaridi 5 Quinidi 6 Sexidi	V. 23 L. 24 M. 25	4 Quaridi 5 Quinidi 6 Sexidi	V. 23 L. 24 M. 25	4 Quaridi 5 Quinidi 6 Sexidi
v. 23 s. 24 D. 25	7 Septidi 8 Octidi 9 Nonidi	V. 24 L. 25 M. 26	7 Septidi 8 Octidi 9 Nonidi	L. 24 M. 25 J. 26	7 Septidi 8 Octidi 9 Nonidi	L. 24 M. 25 J. 26	7 Septidi 8 Octidi 9 Nonidi	V. 24 L. 25 M. 26	7 Septidi 8 Octidi 9 Nonidi	V. 24 L. 25 M. 26	7 Septidi 8 Octidi 9 Nonidi
v. 24 s. 25 D. 26	10 Décidi	L. 26	10 Décidi	M. 26	10 Décidi	M. 26	10 Décidi	V. 25 L. 26	10 Décidi	V. 25 L. 26	10 Décidi
v. 25 s. 26 D. 27	11 Indidi	M. 27	11 Indidi	J. 27	11 Indidi	J. 27	11 Indidi	M. 27	11 Indidi	M. 27	11 Indidi
v. 26 s. 27 D. 28	12 Quindidi	J. 28	12 Quindidi	L. 28	12 Quindidi	L. 28	12 Quindidi	J. 28	12 Quindidi	J. 28	12 Quindidi
v. 27 s. 28 D. 29	13 Sexidi	L. 29	13 Sexidi	M. 29	13 Sexidi	M. 29	13 Sexidi	L. 29	13 Sexidi	L. 29	13 Sexidi
v. 28 s. 29 D. 30	14 Septidi	M. 30	14 Septidi	J. 30	14 Septidi	J. 30	14 Septidi	M. 30	14 Septidi	M. 30	14 Septidi
v. 29 s. 30 D. 31	15 Octidi	J. 31	15 Octidi	L. 31	15 Octidi	L. 31	15 Octidi	J. 31	15 Octidi	J. 31	15 Octidi
v. 30 s. 31 D. 1	16 Nonidi	L. 1	16 Nonidi	M. 1	16 Nonidi	M. 1	16 Nonidi	L. 1	16 Nonidi	L. 1	16 Nonidi
v. 31 s. 1 D. 2	17 Décidi	M. 2	17 Décidi	J. 2	17 Décidi	J. 2	17 Décidi	M. 2	17 Décidi	M. 2	17 Décidi
v. 1 s. 2 D. 3	18 Indidi	J. 3	18 Indidi	L. 3	18 Indidi	L. 3	18 Indidi	J. 3	18 Indidi	J. 3	18 Indidi
v. 2 s. 3 D. 4	19 Quindidi	L. 4	19 Quindidi	M. 4	19 Quindidi	M. 4	19 Quindidi	L. 4	19 Quindidi	L. 4	19 Quindidi
v. 3 s. 4 D. 5	20 Sexidi	M. 5	20 Sexidi	J. 5	20 Sexidi	J. 5	20 Sexidi	M. 5	20 Sexidi	M. 5	20 Sexidi
v. 4 s. 5 D. 6	21 Septidi	L. 5	21 Septidi	L. 5	21 Septidi	L. 5	21 Septidi	J. 5	21 Septidi	J. 5	21 Septidi
v. 5 s. 6 D. 7	22 Octidi	M. 6	22 Octidi	M. 6	22 Octidi	M. 6	22 Octidi	L. 6	22 Octidi	L. 6	22 Octidi
v. 6 s. 7 D. 8	23 Nonidi	J. 6	23 Nonidi	L. 6	23 Nonidi	L. 6	23 Nonidi	M. 6	23 Nonidi	M. 6	23 Nonidi
v. 7 s. 8 D. 9	24 Décidi	L. 7	24 Décidi	M. 7	24 Décidi	M. 7	24 Décidi	J. 6	24 Décidi	J. 6	24 Décidi
v. 8 s. 9 D. 10	25 Indidi	M. 8	25 Indidi	J. 7	25 Indidi	J. 7	25 Indidi	L. 7	25 Indidi	L. 7	25 Indidi
v. 9 s. 10 D. 11	26 Quindidi	L. 8	26 Quindidi	L. 8	26 Quindidi	L. 8	26 Quindidi	M. 8	26 Quindidi	M. 8	26 Quindidi
v. 10 s. 11 D. 12	27 Sexidi	M. 9	27 Sexidi	J. 8	27 Sexidi	J. 8	27 Sexidi	L. 8	27 Sexidi	L. 8	27 Sexidi
v. 11 s. 12 D. 13	28 Septidi	L. 9	28 Septidi	M. 9	28 Septidi	M. 9	28 Septidi	J. 8	28 Septidi	J. 8	28 Septidi
v. 12 s. 13 D. 14	29 Octidi	J. 9	29 Octidi	L. 9	29 Octidi	L. 9	29 Octidi	M. 9	29 Octidi	M. 9	29 Octidi
v. 13 s. 14 D. 15	30 Nonidi	M. 10	30 Nonidi	J. 9	30 Nonidi	J. 9	30 Nonidi	L. 9	30 Nonidi	L. 9	30 Nonidi
v. 14 s. 15 D. 16	31 Décidi	L. 10	31 Décidi	M. 10	31 Décidi	M. 10	31 Décidi	J. 9	31 Décidi	J. 9	31 Décidi

est appelé le Primidi.  
Le jour intermédiaire qui doit terminer cette période  
est appelé le Jour de la Révolution. Il pour cet  
de là. Après les cinq complémentaires.  
XI. Le jour de mensure à mensure se trouve en ce par-  
tir. Le plus petit nombre commun multiple de la durée  
de l'année républicaine est de 12096 jours. Ce nombre se  
divise par le nombre de jours de chaque mois, et le  
résultat est le nombre de fois que le jour de mensure se  
trouve dans l'année républicaine.  
XII. Le nombre de jours de mensure qui se trouvent dans  
une année républicaine est de 12096 jours.  
XIII. Le nombre de jours de mensure qui se trouvent dans  
une année républicaine est de 12096 jours.

est appelé le Primidi.  
Le jour intermédiaire qui doit terminer cette période  
est appelé le Jour de la Révolution. Il pour cet  
de là. Après les cinq complémentaires.  
XI. Le jour de mensure à mensure se trouve en ce par-  
tir. Le plus petit nombre commun multiple de la durée  
de l'année républicaine est de 12096 jours. Ce nombre se  
divise par le nombre de jours de chaque mois, et le  
résultat est le nombre de fois que le jour de mensure se  
trouve dans l'année républicaine.  
XII. Le nombre de jours de mensure qui se trouvent dans  
une année républicaine est de 12096 jours.  
XIII. Le nombre de jours de mensure qui se trouvent dans  
une année républicaine est de 12096 jours.

est appelé le Primidi.  
Le jour intermédiaire qui doit terminer cette période  
est appelé le Jour de la Révolution. Il pour cet  
de là. Après les cinq complémentaires.  
XI. Le jour de mensure à mensure se trouve en ce par-  
tir. Le plus petit nombre commun multiple de la durée  
de l'année républicaine est de 12096 jours. Ce nombre se  
divise par le nombre de jours de chaque mois, et le  
résultat est le nombre de fois que le jour de mensure se  
trouve dans l'année républicaine.  
XII. Le nombre de jours de mensure qui se trouvent dans  
une année républicaine est de 12096 jours.  
XIII. Le nombre de jours de mensure qui se trouvent dans  
une année républicaine est de 12096 jours.

est appelé le Primidi.  
Le jour intermédiaire qui doit terminer cette période  
est appelé le Jour de la Révolution. Il pour cet  
de là. Après les cinq complémentaires.  
XI. Le jour de mensure à mensure se trouve en ce par-  
tir. Le plus petit nombre commun multiple de la durée  
de l'année républicaine est de 12096 jours. Ce nombre se  
divise par le nombre de jours de chaque mois, et le  
résultat est le nombre de fois que le jour de mensure se  
trouve dans l'année républicaine.  
XII. Le nombre de jours de mensure qui se trouvent dans  
une année républicaine est de 12096 jours.  
XIII. Le nombre de jours de mensure qui se trouvent dans  
une année républicaine est de 12096 jours.

# La calendrier Républicain

Converti au format calendrier grégorien.

```
const AN_REGEX = /An (\d+)/i;

export function normalizeYear(year) {
  const m = year.match(AN_REGEX);

  if (!m)
    return +year;

  const nb = m[1];

  if (nb < 2 || nb > 14)
    throw Error(
      `toflit18.republican_calendar.normalizeYear: invalid year ${year}.`
    );

  return 1792 + (+nb);
}
```

# Contrôle de version des données

- `git` - contrôle de version pour les codes sources
- Appliqué à la gestion de corpus de données
- Implique des fichiers texte brut

# fichier texte brut

- Les fichiers texte brut facilitent les traitements informatiques.
- CSV, JSON, XML sont des formats de fichier texte brut.
- Non, XLS, XLSX et ODT n'en sont pas.

TRADE WITH PRINCIPAL COUNTRIES.

PORTUGAL—IMPORTS.

TOTAL VALUE of IMPORTS for HOME CONSUMPTION (MERCHANDISE

COUNTRIES.	1886.	1887.	1888.	1889.	1890.
	Milreis.	Milreis.	Milreis.	Milreis.	Milreis.
Russia	471,000	500,000	887,000	1,085,000	448,000
Norway and Sweden	741,000	876,000	983,000	1,351,000	1,236,000
Germany	4,081,000	4,508,000	4,706,000	5,556,000	6,303,000
Holland	423,000	330,000	309,000	348,000	411,000
Belgium	1,537,000	1,522,000	1,445,000	1,526,000	2,246,000
United Kingdom	12,125,000	12,250,000	12,300,000	13,914,000	13,263,000
France	5,124,000	4,955,000	4,980,000	6,010,000	6,862,000
Switzerland	—	133,000	195,000	267,000	275,000
Spain	2,593,000	2,265,000	2,551,000	3,240,000	2,931,000
Italy	799,000	437,000	675,000	675,000	796,000
Austria	8,000	140,000	371,000	567,000	488,000
Egypt	—	325,000	308,000	216,000	89,000
Morocco	424,000	316,000	239,000	251,000	444,000
China	—	106,000	243,000	263,000	266,000
United States	4,978,000	5,307,000	4,484,000	5,065,000	5,148,000
Brazil	2,008,000	1,874,000	2,148,000	1,803,000	1,946,000

TRADE WITH PRINCIPAL COUNTRIES.

PORTUGAL—IMPORTS.

only) into PORTUGAL, distinguishing PRINCIPAL COUNTRIES.

1891.	1892.	1893.	1894.	COUNTRIES.
Milreis.	Milreis.	Milreis.	Milreis.	
757,000	267,000	408,000	229,000	Russia.
1,103,000	817,000	1,465,000	1,375,000	Norway and Sweden.
5,163,000	2,785,000	4,440,000	4,238,000	Germany.
409,000	328,000	362,000	267,000	Holland.
1,292,000	888,000	1,051,000	1,183,000	Belgium.
11,780,000	9,303,000	10,872,000	9,809,000	United Kingdom.
5,321,000	3,434,000	3,826,000	3,756,000	France.
241,000	221,000	178,000	220,000	Switzerland.
2,402,000	1,701,000	2,837,000	2,930,000	Spain.
754,000	788,000	365,000	551,000	Italy.
469,000	266,000	233,000	226,000	Austria.
311,000	439,000	458,000	430,000	Egypt.
447,000	106,000	213,000	181,000	Morocco.
229,000	228,000	288,000	260,000	China.
5,232,000	6,038,000	7,291,000	5,761,000	United States.
2,054,000	1,911,000	2,428,000	2,484,000	Brazil.

Statistical abstract, P. 166 @ Internet Archive

PORTUGAL—EXPORTS.

TOTAL VALUE of EXPORTS of DOMESTIC PRODUCE (MERCHANDISE

COUNTRIES.	1886.	1887.	1888.	1889.	1890.
	Milreis.	Milreis.	Milreis.	Milreis.	Milreis.
Russia	296,000	318,000	306,000	336,000	323,000
Norway and Sweden	213,000	223,000	299,000	287,000	187,000
Denmark	206,000	160,000	134,000	208,000	231,000
Germany	1,318,000	1,620,000	1,903,000	1,987,000	2,066,000
Holland	206,000	231,000	257,000	278,000	282,000
Belgium	346,000	495,000	376,000	469,000	594,000
United Kingdom	6,707,000	6,764,000	7,828,000	8,228,000	7,993,000
France	9,491,000	4,818,000	5,207,000	3,768,000	1,522,000
Spain	1,156,000	1,210,000	939,000	1,105,000	879,000
Italy	164,000	188,000	168,000	197,000	221,000
United States	647,000	646,000	554,000	509,000	730,000
Brazil	4,575,000	3,686,000	4,195,000	4,260,000	5,181,000
Portuguese Possessions : In Africa	540,000	601,000	879,000	1,007,000	1,111,000
In Asia	17,000	19,000	28,000	22,000	27,000
Other Countries	223,000	260,000	370,000	382,000	281,000
TOTAL	26,108,000	21,239,000	23,443,000	23,343,000	21,538,000
	£ 5,874,000	£ 4,779,000	£ 5,275,000	£ 5,252,000	£ 4,846,000

PORTUGAL—EXPORTS.

only) from PORTUGAL, distinguishing PRINCIPAL COUNTRIES.

1891.	1892.	1893.	1894.	COUNTRIES.
Milreis.	Milreis.	Milreis.	Milreis.	
284,000	406,000	621,000	725,000	Russia.
362,000	173,000	219,000	303,000	Norway and Sweden.
302,000	280,000	269,000	304,000	Denmark.
2,308,000	2,221,000	1,961,000	2,053,000	Germany.
301,000	330,000	396,000	390,000	Holland.
442,000	440,000	632,000	1,048,000	Belgium.
7,483,000	8,719,000	6,553,000	6,602,000	United Kingdom.
1,288,000	1,015,000	906,000	757,000	France.
888,000	1,409,000	1,371,000	2,513,000	Spain.
179,000	105,000	187,000	221,000	Italy.
806,000	900,000	841,000	539,000	United States.
5,274,000	6,782,000	7,155,000	5,994,000	Brazil.
1,280,000	1,624,000	2,061,000	2,187,000	Portuguese Possessions : In Africa.
7,000	26,000	34,000	24,000	In Asia.
174,000	201,000	202,000	164,000	Other Countries.
21,379,000	24,631,000	23,408,000	23,924,000	Milreis } TOTAL.
4,810,000	5,542,000	5,267,000	5,383,000	£ }

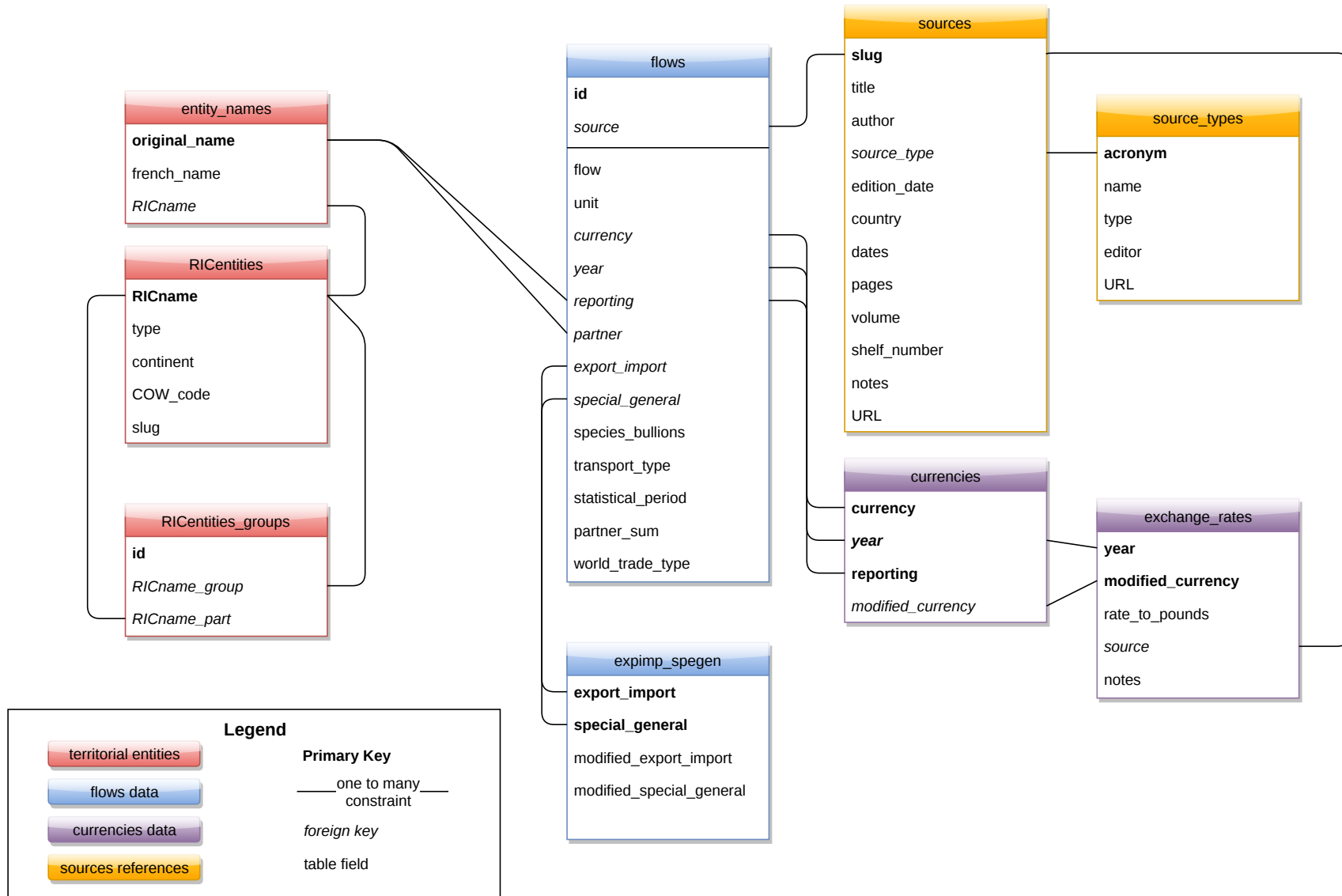
# Prendre soin des données

- Contrôles qualité des données
- Hybridant des approches qualitatives et quantitatives

# Bases de données

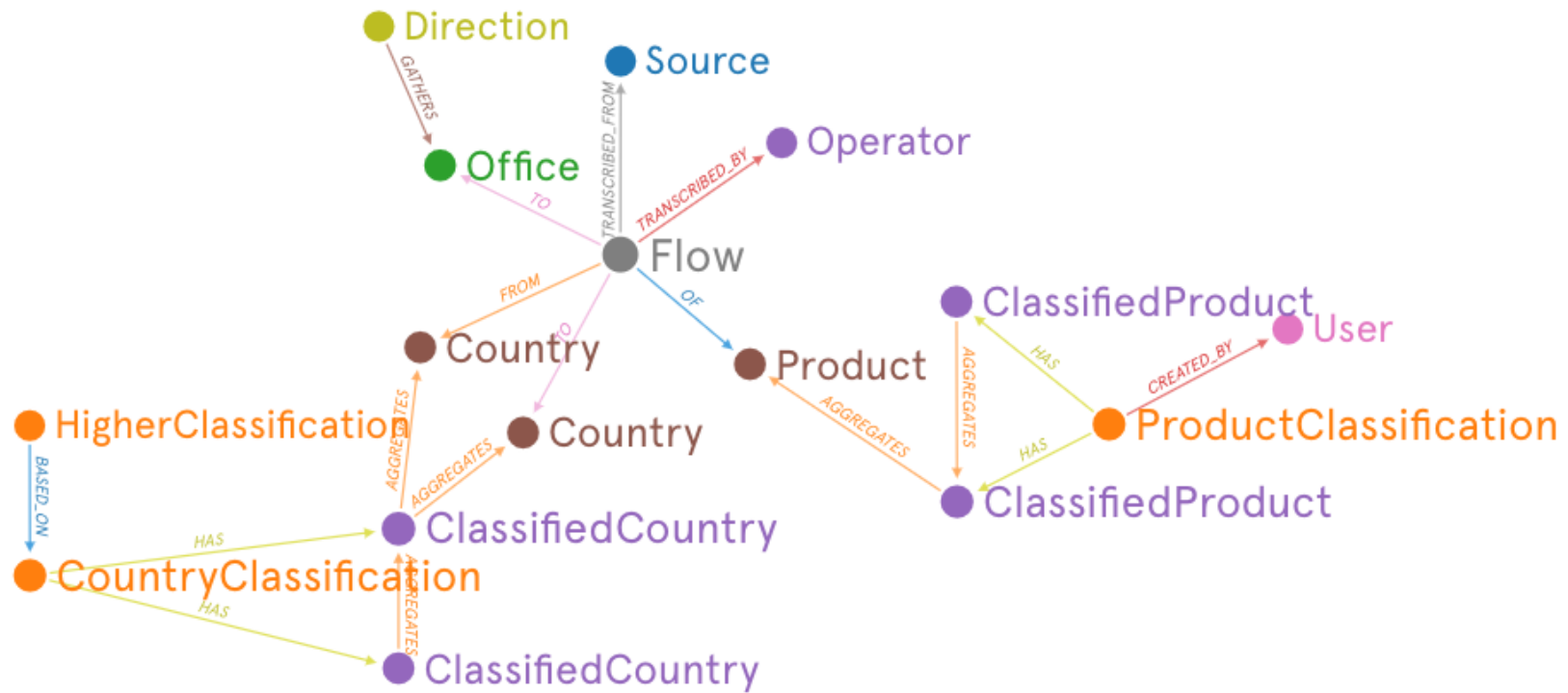
- Les base de données facilitent accès et manipulation
- Le choix de la technologie dépends des usages

# RICardo relational database schema





# Toflit18 NEO4J schema



# Classifications

Un des objectifs principaux de TOFLIT18 :  
créer des classifications dédiées aux questions de  
recherche.

## Classifications

### Products

#### Sources

51 464 items.

**Orthographic Normalization**  
20 608 groups for 51 464 items.  
43 649 / 51 464 (-7 815) classified items (84 %)

**Simplification**  
15 857 groups for 20 608 items.  
2 487 / 20 608 (-21) classified items (99 %)

#### SITC

24 groups for 15 857 items.  
13 365 / 15 857 (-2 492) classified items (84 %)

#### Medicinal products

3 groups for 15 857 items.  
15 826 / 15 857 (-31) classified items (99 %)

### de classification

#### Hamburg classification

39 groups for 15 857 items.

Export

Create From

## Orthographic Normalization



Reset filters

Filter

achia (2 items)

- Achia
- Akia

achia espèce de salade (1 items)

- Akia ; espèce de salade

acide arsenieux (1 items)

- Acide ; arsenieux

acide arsenieux (acide blanc) (2 items)

- Acide arsenieux (Arsenic blanc)
- Acides ; arsenieux (arsenic blanc)

acide benzoïque (1 items)

- Acide ; benzoïque

acide borique (1 items)

- Acide ; borique

# base de données en graphe

Le modèle a été conçu pour proposer des

classifications :

- **hiérarchiques** : agrégation progressive
- **concurrentes** : agrégation dédiée à une question
- **dynamiques** : toute analyse commence par un choix

# la base de données TOFLIT18

- 419729 flux
- 47732 produits
- 843 pays
- 51 bureaux de commerce français
- 120 années
- 807 volumes d'archive

# la base de données RICardo

- 294138 flux
- 1492 RICentities
- 152 années
- 120 monnaies
- 7206 taux de change vers le £
- 73 types de sources (919 volumes)

# L'exploration visuelle au service des données

*Datascape:*

- Des visualisations de données interactives
- Proposant divers points de vue sur les données
- Aide à s'approprier la complexité par les dynamiques d'exploration

Leclercq, C. and Girard, P. (2013). *The Experiments in Art and Technology Datascape*. Collections Électroniques de l'INHA. Actes de Colloques et Livres En Ligne de l'Institut National D'histoire de L'art. INHA <http://inha.revues.org/4926> (accessed 27 October 2015).

# Concevoir un datascape

Ateliers appelés «*data sprints*» avec:

- historiens
- économistes
- développeurs
- designers

Traitant les enjeux de **contenu**, **implémentation** et de **design**

en même temps et au même endroit.

# Analyse Exploratoire de Données

*« The greatest value of a picture is  
when it forces us to notice what we  
never expected to see. »*

*Tukey, J. W.*

Tukey, J. W. (1977). *Exploratory Data Analysis*. Addison-Wesley Publishing Company.



# Visualisation de données

Total trade of reporting **Portugal** from **1846** to **1938**

In this graph, total trade equals the sum of export (import) flows with all partner entities.

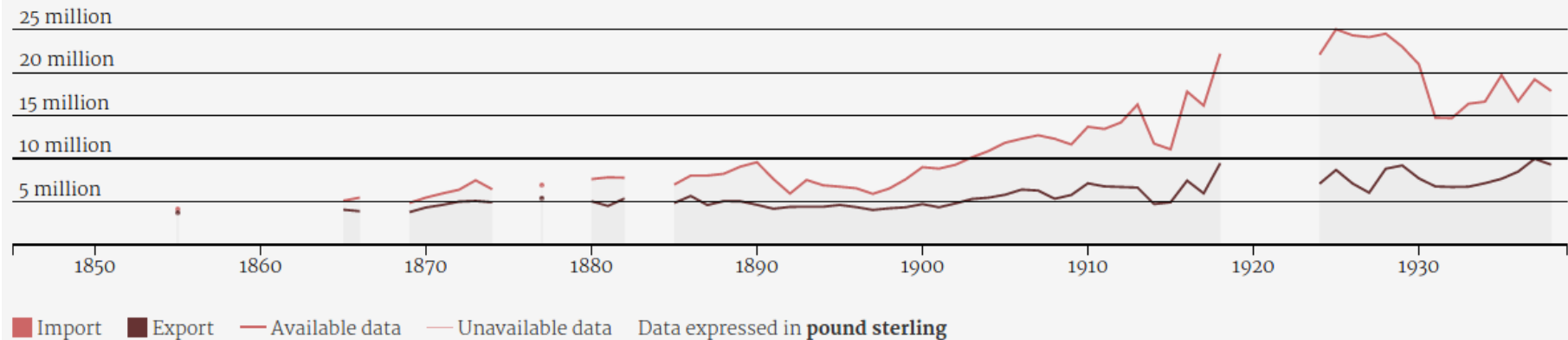


↑ oups !

# Vérifier et corriger les données

Total trade of reporting **Portugal** from **1846** to **1938**

In this graph, total trade equals the sum of export (import) flows with all partner entities.



- corrigé comme indiqué **slide 20** -

# La vue métadonnée de Ricardo

## Reporting entities table

This graph provides detailed information for each reporting entity/year.

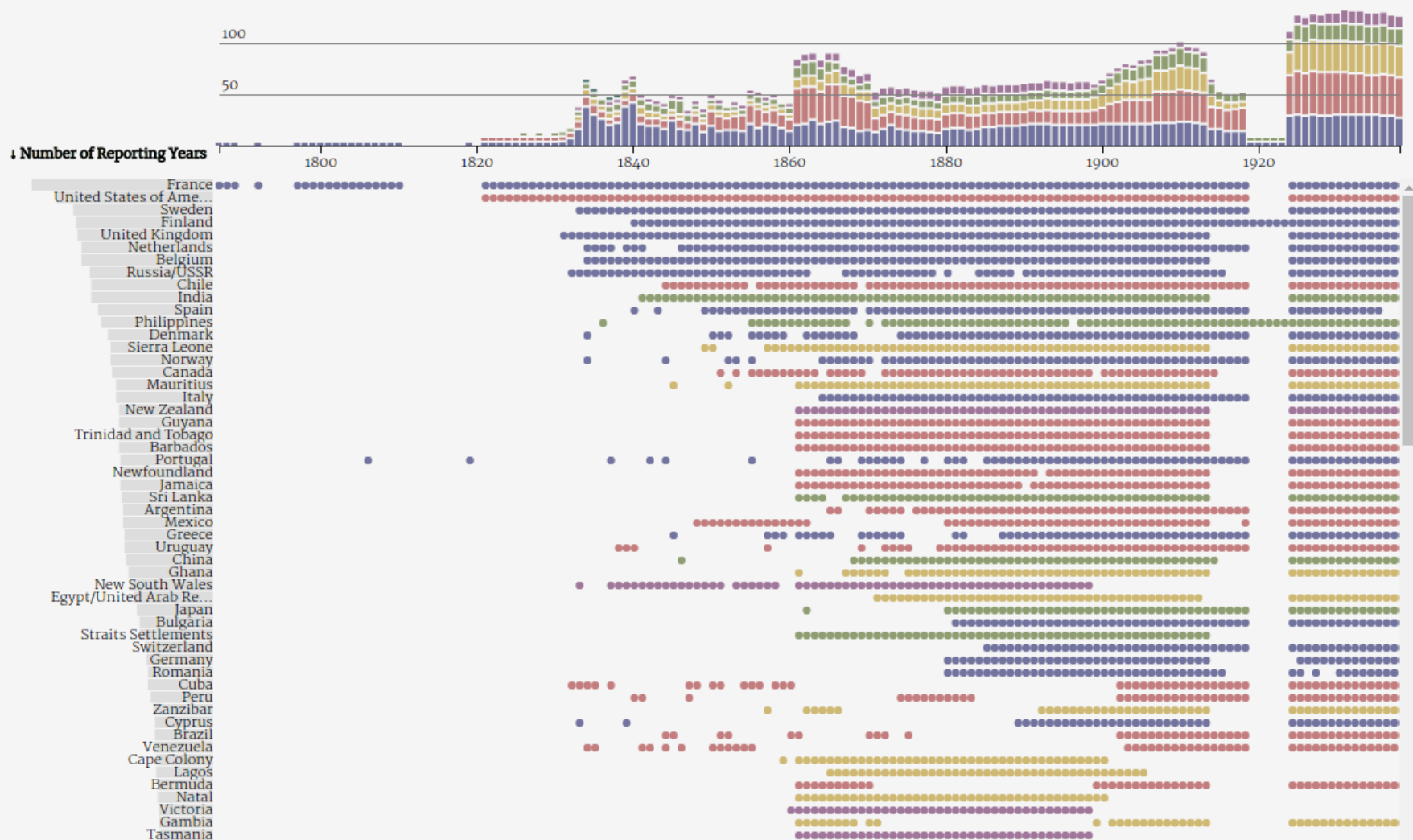
Order by **Number of Reporting Years** Color by **Reporting Continent**

Find Reporting

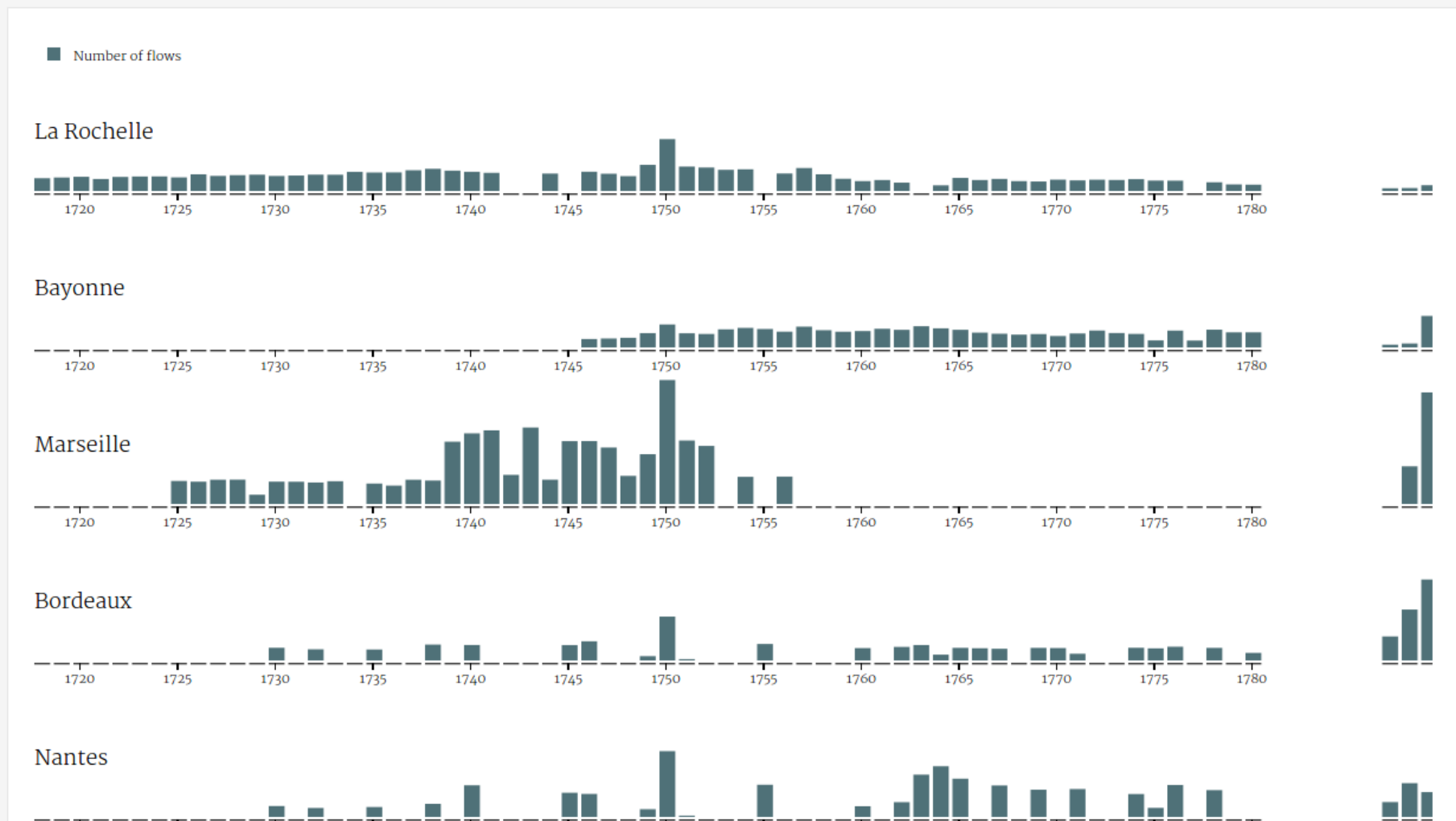
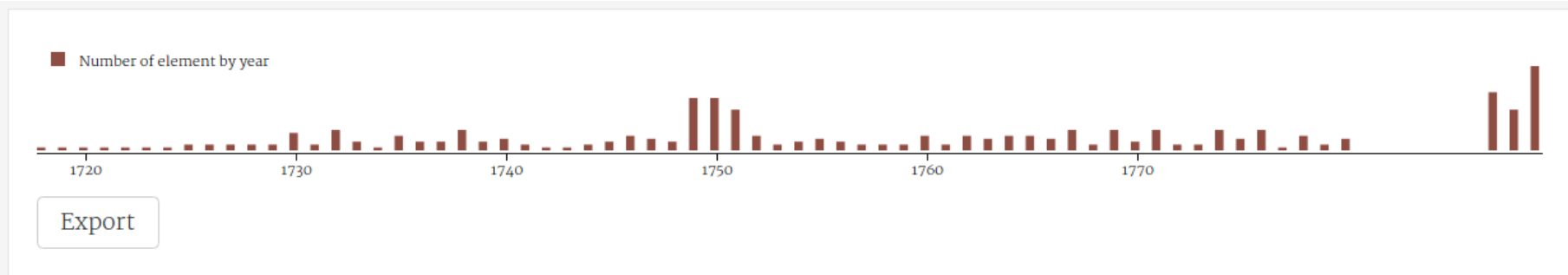
Number of Reportings Color by

**Reporting Continent:** Africa, America, Asia, Europe, Oceania, World. When the 'Reporting type' is a group of countries from different continents, the designated continent is 'World'.

■ Europe ■ America ■ Africa ■ Asia ■ Oceania ■ World



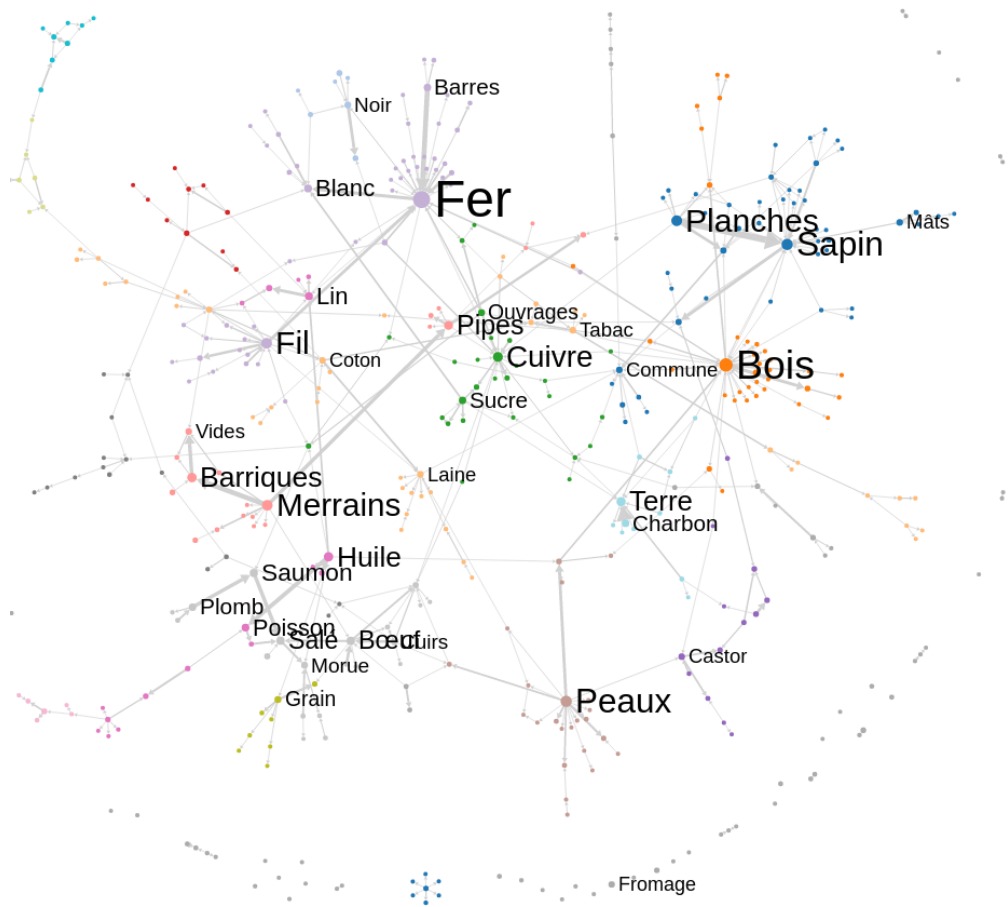
# La vue métadonnée de TOFLIT18



# La chaîne de transformations

volumes d'archive > images > excel > git(csv) > base  
de données > visualisation de données > **csv**

# Un avant goût de l'atelier TOFLIT18/Ricardo



Co-occurrences de termes dans les noms de produits dans les exports de "La Rochelle" entre 1720 et 1729

Rendez vous ce soir à 18h30.

# Les enjeux liés à la visualisation de données historiques

[medialab.github.io/ricardo](https://medialab.github.io/ricardo)

Girard, P., Dedinger, B., Ricci, D., Ooghe-Tabanou, B., Jacomy, M., Plique, G. and Tible, G. (2016). RICardo Project : Exploring XIX Century International Trade. Kraków, Poland [http://ricardo.medialab.sciences-po.fr/Girardetal\\_RICardo\\_dh2016\\_en.pdf](http://ricardo.medialab.sciences-po.fr/Girardetal_RICardo_dh2016_en.pdf).

# Science ouverte

- télécharger les données des visualizations en csv
  - corpus de données seront ouverts en 2017
- CC BY-SA
- RICardo : <http://ricardo.medialab.sciences-po.fr>
  - sources RICardo : [github.com/medialab/ricardo](https://github.com/medialab/ricardo)
  - sources TOFLIT18 : [github.com/medialab/toflit18](https://github.com/medialab/toflit18)
  - cette présentation:  
<http://medialab.github.io/toflit18/ANFmateSHS>



# Bibliographie & Liens (1/2)

Dedinger, Béatrice, and Paul Girard. 2016. *'Exploring Trade Globalization in the Long Run : The RICardo Project'*. Historical Methods.

<http://ricardo.medialab.sciences-po.fr>.

Girard, Paul, Béatrice Dedinger, Donato Ricci, Benjamin Ooghe-Tabanou, Mathieu Jacomy, Guillaume Plique, and Grégory Tible. 2016. *'RICardo Project : Exploring XIX Century International Trade'*. In . Kraków, Poland.

[http://ricardo.medialab.sciences-po.fr/Girardetal\\_RICardo\\_dh2016\\_en.pdf](http://ricardo.medialab.sciences-po.fr/Girardetal_RICardo_dh2016_en.pdf).

# Bibliographie & Liens (2/2)

Latour, Bruno. 1993. *'Le Topofil de Boa-Vista. La Référence Scientifique: Montage Photophilosophique'*. *Raisons Pratiques* 4: 187–216.

Latour, Bruno, Pablo Jensen, Tommaso Venturini, Sébastien Grauwin, and Dominique Boullier. 2012. *“The Whole Is Always Smaller than Its Parts” - a Digital Test of Gabriel Tarde's Monads'*. *The British Journal of Sociology* 63 (4): 590–615. doi:10.1111/j.1468-4446.2012.01428.x.

Tukey, John Wilder. 1977. *Exploratory Data Analysis*. Addison-Wesley Publishing Company.