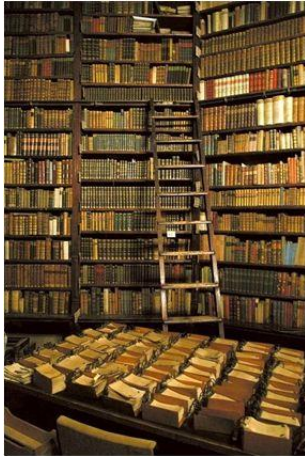


# Construction de bases de données historiques:

## Retours d'expérience



Année	Mois	Jour	Particularité	Montant	Debit	Credit	Balance
1785	Jan	15	...	...	...	...	...
1785	Jan	20	...	...	...	...	...
1785	Jan	25	...	...	...	...	...
1785	Jan	30	...	...	...	...	...
1785	Jan	31	...	...	...	...	...

Lionel Kesztenbaum (INED et PSE)

Gilles Postel-Vinay (PSE)

Ateliers SoDUCo-BNF 4<sup>e</sup> séance

06/11/2023

# Introduction: l'enquête historique, d'hier à aujourd'hui

- De la source (ou des sources) à la (ou aux) base(s) de données
  - Enjeux et défis archivistiques: interprétations et limites des sources, risques téléologiques, etc.
  - Enjeux et défis techniques: difficultés de la collecte, de la constitution et du maintien des bases de données, etc.
  - Enjeux et défis historiographiques: agrégation et normalisation, maintien et perte d'information etc.
- Pérennité et cumulativité
  - Un quatrième enjeu qui a pris dans les dernières décennies une importance croissante c'est la documentation des données et des choix opérés, explicites ou non. Sachant que ces choix portent en eux
    - le risque d'homogénéisation incontrôlée et de perte du lien à la source.
    - Mais qu'ils peuvent aussi permettre, sous certaines conditions, la mise en compatibilité entre les bases, leur transférabilité et partant la cumulativité des recherches.
- Retour d'expérience à de trois exemples de travaux sur l'histoire (quantitative) de la France
  - Numériser et diffuser les enquêtes industrielles du XIXe siècle: enquêtes industrielles.
  - Constituer un échantillon représentatif de la population française sur 150 ans: enquête TRA.
  - Numériser un siècle de recensement français: projet Socface.

# Données historiques numériques - évaluations critiques et perspectives d'utilisations

- Nous sommes invités à nous interroger
  - sur le mode production de données historiques numériques à un moment où ce type de données connaît des changements très rapides. Alors que depuis les années 70, elle reposait surtout sur une saisie classique d'informations extraites de divers corpus, elle prend désormais une autre ampleur en recourant à des processus automatisés susceptibles d'extraction voire d'enrichissement des sources initiales.
  - Et sur les conditions d'appropriation et de réutilisation de ces nouvelles « sources numériques ». Comment et à quelles conditions peut-on s'en emparer efficacement ? Comment les comprendre, évaluer leur « qualité » et leur rapport avec les sources initiales ? Comment comparer ou rendre cumulable les travaux très divers qui peuvent être issus de l'exploitation de ces jeux de données ?

# Le cas d'une enquête des années 1980-90...

- Je vais admettre que les questions que posaient hier la constitution et la circulation de bases de données historiques ne sont pas très différentes de celles qui se posent aujourd'hui ; mieux, que se retourner, fût-ce brièvement, vers elles offre un miroir grossissant aux questions d'aujourd'hui
- Je ne vais pas remonter jusqu'aux époques héroïques des années 1960-70, au temps des bordereaux de dépouillement, des cartes perforées et des grands ordinateurs
- mais pour amorcer la discussion, je vais commencer par évoquer la constitution d'une de ces «sources numériques » réalisée dans les années 1980 sur les enquêtes « industrielles » de la France au XIXe siècle.

## Pourquoi les enquêtes industrielles de la France du XIXe

- J'étudiais à l'époque le fonctionnement des marchés du travail dans la France d'alors, c'est-à-dire dans un monde où il était encore courant de basculer d'un emploi agricole à un emploi industriel au fil des saisons et pour cela j'avais cherché avec plusieurs collègues à construire et à analyser des corpus de données numériques élaborés à partir des enquêtes réalisées par la SGF sur l'agriculture et sur l'industrie.
- S'agissant de l'industrie, en particulier, la tâche paraissait simple puisque, dans les années 1840 puis dans les années 60, la SGF avait réalisé 2 enquêtes qui toutes deux répondaient à un questionnaire très proche et qui l'une et l'autre étaient publiées sous formes de tableaux.
  - Elles visaient un recensement exhaustif du secteur industriel.
  - Elles étaient organisées à l'échelle de l'établissement.
  - et chaque établissements était défini par un ensemble de variables quasi identiques dont une part sur l'emploi et sa saisonnalité

Constituer à partir de la source initiale imprimée la nouvelle «source numérique» paraissait particulièrement simple : transformer les 2 enquêtes en des tableaux croisés portant sur une dizaine de milliers d'établissements définis par une vingtaine de variables.

Dans un 2<sup>e</sup> temps, assortir ce premier travail d'un descriptif assez précis pour permettre dans de bonnes conditions l'appropriation et la réutilisation de ces nouvelles « sources numériques » paraissait aussi assez simple.

- De là un livre de présentation publié en 2000 (Chanut et al. Paris, Ed de l'EHESS, 2000) complété par un dvd puisque c'était à l'époque le seul moyen pour rendre accessible le travail à des tiers...

## **Mais, à l'expérience, constituer la nouvelle « source numérique » et la rendre réappropriable s'est révélé être une opération très loin d'aller de soi (1)**

Les problèmes rencontrés ont été nombreux :

- Ces enquêtes ne sont pas exhaustives (cf. Paris en 1840)(cf. les mines...)
- Elles peinent à définir le champ qu'elles doivent couvrir. En particulier:
  - elles peinent à tracer la frontière entre industrie et artisanat. Cf. C. Lemerancier et l'enquête de la Chambre de Commerce.
  - elles ont eu du mal à décider si elles devaient recenser tous les établissements ou fixer un seuil au-delà d'une certaine taille et, si oui, quel seuil (>1 ouvrier, ou à 10 ? Ou à 20 ?).
- Elles ont hésité sur le nombre de questions pertinentes.
- Enfin –et c'est peut-être plus surprenant – elles ont hésité sur la frontière entre « industrie » et « agriculture », une frontière très poreuse quand il s'agissait de l'emploi car, à l'époque, , comme je l'ai dit, une partie souvent importante de la main d'œuvre bascule d'un secteur à l'autre au rythme des saisons (au moment des récoltes).

## **Mais, à l'expérience constituer la nouvelle « source numérique » et la rendre réappropriable s'est révélé être une opération très loin d'aller de soi (2)**

- Parmi les problèmes rencontrés certains ont été assez simples à documenter.
  - Les seuils.
  - Les frontières internes du secteur de l'industrie (définition des branches) etc...
- D'autres ont soulevé plus de nombreux problèmes. Pour les comprendre il faut se souvenir de ce que Jean-Claude Perrot appelait de ses vœux : une histoire matérielle de l'abstraction.



# Une histoire matérielle de l'abstraction (1)

- Parenthèse : Lucien Febvre 1932 « **Les textes ... mais tous les textes ...** »
- « Les textes » : jusqu'ici on en a donné une définition restrictive : on est parti des volumes imprimés.
- « Mais tous les textes »... Et, pour cela, remonter de l'imprimé à la source primaire et d'abord aux questionnaires (quitte à aller aux AD... où ils ne sont pas tjs conservés). Pb :
- Certaines questions posées par les questionnaires pouvaient être immédiatement transcrites dans les colonnes des tableaux publiés. Dans ce cas tout était simple.
- Mais d'autres réponses ont dû être ré-élaborées en bureau par les responsables des enquêtes (ex: les branches)... selon des modalités qui n'allaient pas de soi...
- Enfin, les questionnaires comportaient plus de questions que celles reprises dans les publications (par exemple les questions sur l'origine des matières premières et la destination des produits. Le passage à l'imprimé a donc imposé des choix à la fois intellectuels et matériels : quelles variables retenir dans l'imprimé compte tenu de l'espace disponible (une page du volume comporte un nombre fini de colonnes).

## Une histoire matérielle de l'abstraction (2) Des questions initiales à la source publiée

				Renseignement présent		Renseignement obtenu par	
				au questionnaire		agrégation en bureau	
				1840	1860	1840	1860
Département, Arrondissement				OUI	NON	non sauf exceptions	OUI
Commune				OUI	OUI	non sauf exceptions	NON
Nom du fabricant				OUI	NON	non sauf exceptions	NON
branche d'activité				NON	NON	OUI	OUI
Valeur locative et montant de la patente				OUI	OUI	non sauf exceptions	OUI
Matières premières		nature/quantités/valeur		valeur seulement	OUI	non sauf exceptions	OUI
Produit		nature/quantités/valeur		valeur seulement	OUI	non sauf exceptions	OUI
Débouché ou lieu de destination des produits				OUI	NON	NON	NON
Mois chômés				NON	OUI	NON	OUI
Ouvriers		H/F/E nb et salaire/jour		OUI	OUI	non sauf exceptions	OUI
Moteurs		eau, vent, animaux, vapeur		OUI	OUI	non sauf exceptions	OUI
Feux				OUI	OUI	non sauf exceptions	OUI
Machines				OUI	OUI	non sauf exceptions	OUI

# Bulletin individuel de la 1<sup>e</sup> enquête

*Arrondissement*  
de Melun.

*Statistique de France — Industrie*

*Commune*  
de Fontaine-Combaud.

Valeur locative de l'établissement  
250 francs

*Département de Seine-et-Marne.*

Montant de la patente (Calculé en la valeur de l'établissement)

*Fabrique et Manufacture de toiles, linges et corraons.*

Nom de fabricant, entrepreneur ou manufacturier  
M. Bourgeois Jean Louis Viret.

Production industrielle	Nature des objets	Leur quantité en unités ou kilogrammes ou mètres ou hectolitres	Valeur de chaque unité en francs et centimes	Valeur totale, par approximation, en francs	Lieux Usages de matières premières
Matières premières employées	Crus blancs	150 mètres cubes	2 francs la pièce cube	300 <sup>f</sup>	Des la laine de la Fabrique et des la commune.
	Crus jaunes	250 mètres cubes	1 <sup>f</sup> 50 <sup>c</sup>	375 <sup>f</sup>	
	Sable	30 mètres cubes	3 <sup>f</sup> 00 <sup>c</sup>	90 <sup>f</sup>	
Produits fabriqués manufacturés	Toiles	150000	28 francs la pièce	4200 <sup>f</sup>	Fabriqués en laine de fabrication de produits dans les communes situées de la fabrique en sur Paris
	Draps	75000	40 <sup>f</sup>	3000 <sup>f</sup>	
	Corraons	100000	28 <sup>f</sup>	2800 <sup>f</sup>	

*Recevoir*

Les unités	}	Crus blancs	150
		Crus jaunes	250
		Sable	30
Les valeurs premières en francs et centimes	}	Crus blancs	300 <sup>f</sup>
		Crus jaunes	375 <sup>f</sup>
		Sable	90 <sup>f</sup>

14 sur l'annuaire N. 18 1888

*Matières*

Matières	}	Crus blancs	150
		Crus jaunes	250
		Sable	30
		Crus blancs	150000
		Draps	75000
		Corraons	100000

15) Les francs en unités de devises

*Four*

Four	}	Crus blancs	150
		Crus jaunes	250
		Sable	30

17) 9303

*Machines*

Machines	}	Crus blancs	150
		Crus jaunes	250

à l'usage de la fabrique N. 1888  
D'après l'annuaire N. 1888  
M. Bourgeois



## Une histoire matérielle de l'abstraction (3)

- Exemple 1/ Quand ont été publiés les résultats de la première enquête, une page imprimée ne pouvait contenir toutes les variables et il a fallu choisir entre publier le nom des propriétaires et leur adresse (deux colonnes) **ou** publier les données sur la localisation des consommations intermédiaires et des produits :
  - La décision a été prise de renoncer aux secondes pour publier l'identité et le domicile des propriétaire, puisque la logique d'ensemble reposait sur le choix de publier les information établissement par établissement.
  - Conséquences :
    - Certaines positives : on observe ainsi non seulement les établissements mais aussi les entreprises (puisque l'enquête identifie son propriétaire) et son statut juridique.
    - D'autres négatives : ont ainsi surgi de nombreux conflits avec les enquêtés (très hostiles à ce que les renseignements recueillis soient rendus publics).
    - Si nombreux que c'est la décision inverse qui est prise en 1860 quitte à changer la présentation des résultats et ne les publier qu'agrégés pour les établissements d'une même branche dans un même arrondissement. Ce faisant, on perd notamment la distinction entre établissement et entreprise.

## Une histoire matérielle de l'abstraction (4)

- Exemple 2/ ces deux enquêtes prennent la suite d'opérations antérieures plus ou moins abouties mais avec un questionnaire voisin.
- Elles cherchaient notamment à cerner les flux saisonniers de la main d'œuvre entre industrie et l'agriculture.
  - La 1e enquête industrielle a minoré ce point et décidé de privilégier l'information sur le nom et le domicile des propriétaires chois qui, on vient de le dire, structure la publication.
  - La 2e enquête ayant renoncé à individualiser les propriétaires des établissements, l'espace laissé libre dans la page de la publication a pu être employé à publier l'information sur l'activité saisonnière et donc les mois où est employé la main d'œuvre. De là de nombreuses difficultés pour agréger à l'échelle de chaque branche les résultats de chaque établissement tels qu'ils figuraient au questionnaire.
  - Pb intéressant, mais demandant de tels développements que j'ai fini par renoncer à les expliciter dans le livre sus mentionné...

## Bilan et perspectives (1)

- Le premier avantage de ces nouvelles « sources numériques » a été de rendre accessibles les info micros sur les établissements (ou au moins sur les établissements d'une même branche dans un même arrondissement). C'était un saut qualitatif par rapport aux recherches antérieures qui n'avaient retenu de ces enquêtes que leurs résultats agrégés à l'échelle du pays pour établir les comptabilités nationales rétrospectives, ce qui était l'un des grands objectifs de toute une génération à l'époque.
- L'approche s'est ensuite révélée fructueuse quand les recherches ont porté sur les établissements de branches les mieux cernées par les enquêtes (le textile cotonnier en particulier: cf. les articles de Pierre Sicsic et, récemment, ceux de Reka Juhasz (reliant pour ce type d'établissements les enquêtes qui se sont succédées de l'Empire à celle de 1840)(effet durables du blocus continental; quel rythme d'adoption de techniques nouvelles).
- Mais on peut imaginer bien d'autres objectifs. Pour cela, à nouveau revenons à Lucien Febvre. « **Les textes** » dit-il (ici les enquêtes publiées)...
- puis il précise « **Les textes ...mais tous les textes** » (ici les questionnaires originaux)...
- Mais il ajoute :

## Bilan et perspectives (2)

- **« Les textes, sans doute : mais tous les textes... Les textes, évidemment : mais pas rien que les textes. Les documents aussi, quelle qu'en soit la nature » ; et Febvre de citer comme exemple les documents qu'utilise ... la linguistique ... ou la palynologie « en quête de pollen millénaire ». « L'histoire fait son miel avec lui. L'histoire qui s'édifie, sans exclusion, avec tout ce que l'ingéniosité des hommes peut inventer et combiner pour suppléer au silence des textes, aux ravages de l'oubli. ».**
  - Febvre en 1932 pense à déplacer les interrogations des historiens en recourant à des croisements entre des sources classiques et des sources nouvelles : en 1967 avec ses recherches sur le climat Le Roy Ladurie (1967) se saisit de ce programme mais aujourd'hui d'autres enquêtes deviennent possibles, en particulier en tirant parti des informations non ambiguës des enquêtes.
  - Cf. l'enquête de 1840 dont les informations localisées (commune par commune) peuvent être réanalysées en les croisant avec les sources qui désormais renseignent sur les conditions climatiques locales, sur la présence d'une rivière ou d'un fleuve (plus ou moins puissance) et donc l'existence de ressources hydrauliques en un temps où c'est encore l'énergie dominante...etc.



# L'enquête TRA: une enquête, des données, un programme de recherche

- Une base de données historiques qui a 40 ans d'histoire
  - Le projet initial: l'enquête « 3000 familles », portée par Jacques Dupâquier au début des années 80.
  - Différentes extensions, en premier lieu autour des données fiscales (Denis Kessler).
- Un contexte de développement rapide de bases de données historiques individuelles
  - Projets initiaux de la démographie historique: Louis Henry, Princeton European Fertility Project, The Cambridge Group for the History of Population and Social Structure, etc.
  - Grandes bases de données individuelles: Balsac et PRDH, Scania database, Eurasian project, etc.
- Changement d'échelle et de finalité
  - Vers une logique cumulative : l'enquête TRA (appariement avec les données militaires, recensement, retraités de la fonction publique etc.).
  - Une base de données individuelles représentative.
  - Avec le temps, la base a changé d'objectifs : **une matrice.**

# L'enquête TRA : une base individuelle et nationale (1)

- Tous les individus dont le noms de famille commence par les lettres TRA.
- Deux sources principales
  - L'état civil et surtout les mariages : Jacques Dupâquier puis Jean-Pierre Pélissier et Danielle Rébaudo à l' EHESS, CNRS, INRA et à l'INED – “l'enquête des 3000 familles” 1984-1992.
  - Les sources fiscales (l'Enregistrement : taxation des successions) : Denis Kessler au CNRS et à l'EHESS, 1985-1992, puis Gilles Postel-Vinay et al. at INRA-EHESS-INED, depuis 1995.
- Un échantillon représentatif de la population française :

échelle  $\approx$  1 pour 1100

# L'enquête TRA : une base individuelle et nationale (2)

## ➤ Pourquoi TRA?

- Représentatif à l'échelle nationale.
- Stable dans le temps.
- Facile à lire.
- Effectifs suffisants, mais pas excessifs.

## ➤ Collecte des données

- Tous les mariages de 1802-1902 – et un sous-échantillon entre 1903-1993.
- Sources fiscales : tous les morts entre 1798-1939 et un sous-échantillon 1940-1959.

## ➤ Reconstitutions de familles.

- “3000 familles” : l'objectif initial à partir de 3000 couples formés entre 1800 et 1830 choisis au hasard
- Reconstitutions en lignées masculines descendantes.

# L'enquête TRA : une base individuelle et nationale (3)

## ➤ Un double objectif

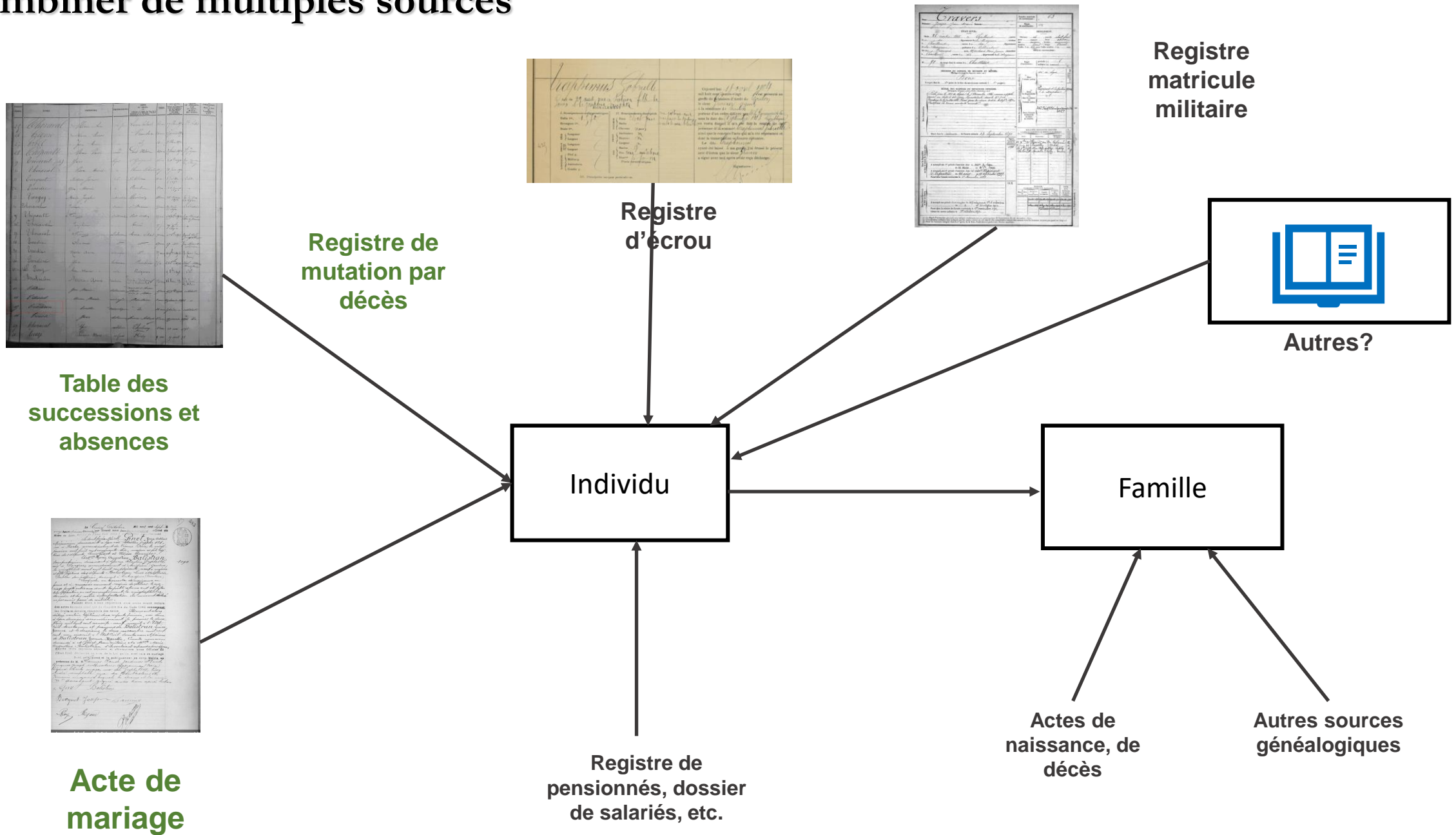
- Trajectoire de vie individuelle (naissance – mariage – décès)
- Reconstitutions de familles sur 7 générations.

## ➤ De la démographie à l'ensemble de l'histoire économique et sociale

- Transformations de la société française dans la longue durée.
- Industrialisation, urbanisation, vieillissement.
- Inégalités, mobilité sociale.

## ➤ Changement de logique : appariement avec d'autres sources nominatives, une logique cumulative.

# Combiner de multiples sources



# Un projet collectif

- Construction de la base de données démographiques : Jacques Dupâquier, Jean-Pierre Pélissier, Myriam Provence..... et des centaines de bénévoles (correspondants de l'enquête).
- Construction de la base de données fiscales : Denis Kessler, Eliane Monge.
- Élargir la reconstitution de la famille : Pascal Cristofoli, Clément Dherbécourt, Jean-Brieu Delbos, Viviani Huerta, Danièle Rébaudo, Eric Zolotukine.
- Élargir l'ensemble des données sur la richesse : Linda Amiri, Guillaume Bazot, Charlotte Coutant, Sarah Favre, Alena Lapatniova.
- La reproduction familiale : Noël Bonneuil, Michel Dürr.
- Les migrations dans la durée : Didier Blanchet, Paul-André Rosental.
- La mobilité sociale à long terme : Marco Van Leeuwen, Ineke Maas.
- Micro-étude de la transmission de la richesse : Luc Arrondel, Cyril Grange.
- Evolution de la structure de la richesse et des professions : Jérôme Bourdieu, Lionel Kesztenbaum, Martha Ménendez, Gilles Postel-Vinay, Akiko Suwa-Eisenman.



# Bilan et perspectives

- Addition de nouvelles sources à partir de la structure nominative: une « enquête » construite pour être diffusée
  - Un socle commun qui ouvre de nombreuses possibilités d'extension (mais les limite aussi: il faut que la source soit (aisément) accessibles nominativement).
  - Et en même temps un élargissement qui implique une appropriation de l'existant, parfois complexe, souvent long.
  
- Suivre les individus dans le temps et l'espace
  - Extension (et généralisation) du principe initial de l'enquête.
  - Reconstitution de trajectoires individuelles et familiales.
  - Ouverture à de nombreuses thématiques.
  
- Mais aussi nombreux défis: instabilités des noms (et prénoms), données manquantes, stabilité des sources et du vocabulaire...

# Retranscrire les listes individuelles du recensement : le projet Socface

- Le projet Socface: collecter, traiter, retranscrire, organiser et analyser l'ensemble des listes nominatives du recensement de 1836 à 1936 (20 recensements).
- Socface produira une base de données des individus ayant vécu en France entre 1836 et 1936.
  - ❖ Construire une histoire économique et sociale à la fois individuelle et nationale.
  - ❖ Grande souplesse des approches: sur une seule année/coupe, plusieurs années combinées, avec appariement des individus entre les recensements ou non...
- Un projet de taille?
  - ❖ 20 millions d'images à traiter + décentralisation (une centaine de dépôt d'archives).
  - ❖ Un projet rendu possible par les méthodes automatisés.
  - ❖ Mais ces méthodes modifient aussi la substance du projet: façon de travailler, rendu en masse...



DESIGNATION		NUMÉROS. PAR QUARTIER, VILLAGE, LIGNES ou RUE.			NOMS	PRÉNOMS.	ANNÉE	LIEU	NATIONALITÉ.	SITUATION PAR RAPPORT	PROFESSION.	REMARQUES
des usines, villages ou hameaux	des rues, des villages, des hameaux	des numéros	des ménages	des individus	DE FAMILLE.		de naissance.	de naissance.	LITTÉ.	ou chef de ménage.		13
				2	Eribout	Lucie	1886	Reims	Fr	épouse	ans	
				3	Eribout	Henri	1909	Paris	"	enfant	"	
				4	Eribout	Alphonse	1912	"	"	"	"	
			11	5	Eribout	René	1916	Aubervilliers	"	"	"	
				6	Eribout	Audré	1920	"	"	"	"	
				7	Pierre	Henriette	1881	Reims	"	bell. soeur	Magasinier	Boulonnais Valenciennes
				8	Pierre	Liliane	1912	Baguillet	"	filie	sans	
				9	Pierre	Lucien	1918	Woisylès	"	neveu	sans	
				1	Parmenier	Justave	1877	Weschateau	"	chef	maçon	
				2	Parmenier	Lucie	1867	Lambach	"	époux	sans	
			12	3	Parmenier	Maximilien	1888	Paris	"	enfant	Confectionnier	Leulman
				4	Parmenier	Emma	1900	W. Denis	"	"	empl.	Minist. Pensions
				5	Parmenier	Ludovic	1902	"	"	"	filie	Alix
				1	Caliez	Jorges	1878	Lill	"	chef	Condamné	Magist.
				2	Caliez	Celine	1884	Stades	"	époux	sans	
				1	Renault	Jean	1869	Bombigny	"	chef	chef équipe	Magist.

Rue des Fillettes

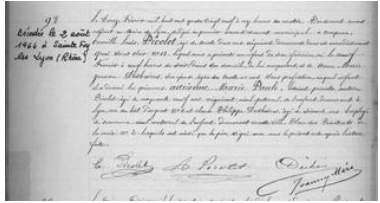
Une liste nominative: Aubervilliers, 1921

# Retranscrire les listes individuelles du recensement : le projet Socface

- Le projet Socface: collecter, traiter, retranscrire, organiser et analyser l'ensemble des listes nominatives du recensement de 1836 à 1936 (20 recensements).
- Socface produira une base de données des individus ayant vécu en France entre 1836 et 1936.
  - ❖ Construire une histoire économique et sociale à la fois individuelle et nationale.
  - ❖ Grande souplesse des approches: sur une seule année/coupe, plusieurs années combinées, avec appariement des individus entre les recensements ou non...
- Un projet de taille?
  - ❖ 20 millions d'images à traiter + décentralisation (une centaine de dépôt d'archives).
  - ❖ Un projet rendu possible par les méthodes automatisés.
  - ❖ Mais ces méthodes modifient aussi la substance du projet: façon de travailler, rendu en masse...

# Le projet Socface dans un temps long

- Le projet vise à être le point de départ des études quantitatives sur la France avant la seconde guerre mondiale.
- Peut aussi permettre de faire le lien avec les données contemporaines (seconde moitié du XXe siècle).
- Réflexion sur les conditions de réutilisation des données
  - Penser en amont la diffusion, techniquement et scientifiquement.
  - Transparence des choix qui sont fait + en faire le moins possible: tentation de « corriger les sources ».
  - Documenter au maximum les étapes de la constitution de la base.
- Développer des indicateurs, outils, marqueurs (...) pour faciliter la réutilisation
  - Indices de « qualité » aisément mobilisables par les sciences sociales.
  - Documenter l'idiosyncrasie de la source: rendre possible comparaison et réutilisation sans les préempter.



Acte de naissance



Registre matricule militaire



Acte de mariage



Registre d'écrou



Registre du commerce et des sociétés



Registre de pensionnés



Acte de décès

Table des successions et absences

Registre de mutation par décès



Naissance

mariage

Décès

Liste nominative

Liste nominative

Liste nominative

Liste nominative

Parents  
Fratrie

Conjoint  
Enfants

# Le projet Socface dans un temps long

- Le projet vise à être le point de départ des études quantitatives sur la France avant la seconde guerre mondiale.
- Peut aussi permettre de faire le lien avec les données contemporaines (seconde moitié du XXe siècle).
- Réflexion sur les conditions de réutilisation des données
  - Penser en amont la diffusion, techniquement et scientifiquement.
  - Transparence des choix qui sont fait + en faire le moins possible: tentation de « corriger les sources ».
  - Documenter au maximum les étapes de la constitution de la base.
- Développer des indicateurs, outils, marqueurs (...) pour faciliter la réutilisation
  - Indices de « qualité » aisément mobilisables par les sciences sociales.
  - Documenter l'idiosyncrasie de la source: rendre possible comparaison et réutilisation sans les préempter.

# Conclusion (1) de la question de recherche à l'enquête, et retour

## ❑ Diversité des sources et des « enquêtes »

- ▲ Enquête industrielle: retracer le cheminement de l'enquête.
- ▲ Enquête TRA: comprendre les apports et limites de chaque source.
- ▲ Projet Socface: gérer la masse et la diversité d'une source unique.

## ❑ Mais des questionnements mutuels

- ▲ Quelle partie du monde social est dévoilée par l'enquête et comment percevoir au mieux extension et limites de l'observatoire.
- ▲ A quelle étape faire la consolidation des informations, le codage des données, la sélection des variables, etc.
- ▲ Linking entre les sources: objectifs et limite, sélection(s) qui en résulte.



# Conclusion (2) Complémentarités et cumulativité

## ❑ Complémentarités à différentes échelles

- ▲ Transmission et réutilisation de l'enquête.
- ▲ Enquêtes individuelle, qualitatives et quantitatives.
- ▲ Données contextuelles: caractéristiques des communes (population, position géographique)) et leurs évolutions (construction d'une école, d'une gare, d'une usine...).

## ❑ Effet de la reconnaissance automatique pour le développement de larges bases de données (historiques)

- ▲ Une forme d'arbitrage entre des données riches (denses) et des données vastes.
- ▲ Question de la représentation: qui est présent? qui est absent? pourquoi?
- ▲ Variation dans la distance entre la source et la base de données... entre la base et les objets.
- ▲ Déconnexion entre la production et la valorisation des données.