

Le data paper, de la littérature grise à l'article de données. Retour sur l'expérience du CDSP et de l'accompagnement des chercheurs à Sciences Po

Anna Egea, Guillaume Garcia

▶ To cite this version:

Anna Egea, Guillaume Garcia. Le data paper, de la littérature grise à l'article de données. Retour sur l'expérience du CDSP et de l'accompagnement des chercheurs à Sciences Po. Un data journal interdisciplinaire pour les sciences humaines et sociales. Enjeux scientifiques et mise en œuvre pratique, Université de Nancy; MSH Lorraine, Mar 2023, Nancy (FR), France. hal-04050447

HAL Id: hal-04050447 https://sciencespo.hal.science/hal-04050447

Submitted on 29 Mar 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Le data paper, de la littérature grise à l'article de données

Retour sur l'expérience du CDSP et de l'accompagnement des chercheurs à Sciences Po

Anna Egea, CSO, CNRS Sciences Po Guillaume Garcia, CDSP, Sciences Po

Journée d'étude "Un data journal interdisciplinaire pour les sciences humaines et sociales - Enjeux scientifiques et mise en œuvre pratique"

Université de Nancy, 10 mars 2023

Plan de l'intervention

Partie 1 : Retour sur l'expérience du CDSP avant l'émergence des data papers

Partie 2 : Retour sur la mise en place d'un service d'accompagnement aux data papers à Sciences Po

Conclusion : Quels enjeux autour des data papers ?

Partie 1

Retour sur l'expérience du CDSP avant l'émergence des data papers

Le contexte originel : la diffusion des enquêtes à Sciences Po

Depuis 2006 au CDSP - essentiellement via des ingénieurs BAP D en documentation

Au départ : données "quantitatives"

A partir de 2013, intégration des données "qualitatives"

Des modèles de documentation mis en oeuvre avant l'émergence des data papers

La banque de données du CDSP sur Datasciencespo SciencesPo

Recherche ▼ Guide d'utilisation Support Français ▼ Se connecter

données qualitatives

(Sciences Po, Centre de données socio-politiques (CDSP), CNRS)

data.sciencespo > Banque de données du CDSP > données qualitatives >

Panel de Caen

Version 3.1



Bidart, Claire; Degenne, Alain; Kornig, Cathel; Lavenu, Daniel; Mounier, Lise; Pellissier, Anne; Le Gall, Didier; Beynier, Dominique; Lemarchand, Clotilde, 2023, "Panel de Caen", https://doi.org/10.21410/7E4/OCUEIL. data.sciencesoo. V3

Citer l'ensemble de données • Pour en apprendre davantage sur le sujet, consulter le document Data Citation Standards [en].

Modalités d'accès à l'ensemble de données →

Contacter le propriétaire Partager

Statistiques d'utilisation sur l' 'ensemble de données (2)

3 téléchargements @

Description @

Élaborée par un groupe de chercheurs coordonnés par Claire Bidart, cette enquête longitudinale par panel étudie les parcours de vie et la dynamique des réseaux sociaux, ainsi que leurs interdépendances, à un moment particulier de la socialisation, au long des transitions vers la vie adulte. Elle porte sur un groupe de jeunes originaires de la région de Caen en Normandie, composé initialement de 87 personnes âgées de 17 à 23 ans, pour moitié filles et pour moitié garçons, alors en classe de terminaie de bac « économique et social » ou de bac professionnel, ou en stage d'insertion. Sur les 20 années qu'a duré l'enquête, le panel est interrogé une première fois en 1995 puis réinterrogé à 5 reprises : en 1998, 2001 et 2004, successivement 73, 56 et 60 de ses jeunes ont à nouveau participé à fenquête. La Sème vague.

Lire la suite de Description [+]

Sujet @

Sciences sociales

Mot-clé ❷

RESEAUX SOCIAUX, SOCIABILITE, JEUNESSE, EVENEMENTS DE VIE

Publication liée @

Claire Bidart, Alain Degenne, Michel Grossetti, La vie en réseau. Dynamique des relations sociales, Paris, PUF, coll. «Le lien social », 2011, 356 p. isbn: 9782130590644

Apports et limites de la documentation "classique" des données

L'enquête, en tant que processus de recherche global, est documentée avec une notice, via des standards de métadonnées de type DDI, ELSST, etc.

Les métadonnées d'une notice d'enquête

SciencesPo

Recherche + Guide d'utilisation Support Français + Se connecter

Métadonnées géospatiales ^

Couverture géographique @

France (la), Normandie, Caen

Unité géographique @

Région

Métadonnées sur les sciences sociales et les sciences humaines ^

Unité d'analyse @

Individu

Univers Q

Panel de 87 jeunes de la région de Caen regroupant en 1995 : des lycéens de section ES, âgés de 17 à 20 ans ; des lycéens

en bac. professionnel, âgés de18 à 23 ans ; des stagiaires, âgés entre 18 et 24 ans .

Méthode temporelle @

Longitudinale: cohorte / basée sur un événement

Responsable de la collecte de

données @

Bidart, Claire

Fréquence @

Tous les trois ans pour les cinq premières vagues, puis huit ans plus tard pour la dernière.

Taille de l'échantillon cible Q

Entretien en face-à-face

Type d'instrument de recherche @

Ouestionnaire semi-structuré

Caractéristiques de la collecte de

données @

Enquête réalisée sur 6 vagues, entre 1995 et 2015, à travers des questionnaires et des entretiens auprès des 87 enquêtés. L'étude se déroulant sur le temps long, l'enquête subie des attritions et le nombre d'enquêté répondant varie selon les

vagues. Les élèves ont été sensibilisés au projet en classe, et devaient ensuite se porter volontaires pour participer au panel. Les volontaires étaient contactés par téléphone afin de procéder au premier entretien. Les élèves étaient recrutés

au sein de trois formations différentes : bac ES, bac professionnel, et stagiaires travaillant déjà.

Mesures visant à minimiser les pertes

Des nouvelles régulières étaient échangées avec les panélistes afin d'éviter l'attrition.

Pondération @ Il n'existe aucune pondération pour ces données qui n'ambitionnent pas d'être représentatives.

Remarques @ Nombre de documents : 926 fichiers numériques ou numérisés. Les documents n'ayant pas été numérisés restent consultables aux archives MSH de Caen Transcription : Documents textuels. Langage documentation : Français.

Anonymisation : Anonymisation réalisée par le chercheur déposant et l'équipe du CDSP.

Apports et limites de la documentation "classique" des données

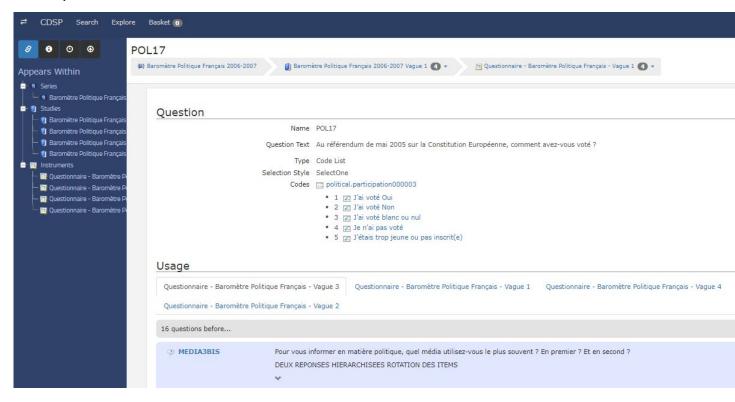
L'enquête en tant que processus de recherche global est documentée via des standards de métadonnées de type DDI, ELSST, etc.

Les descriptions détaillées des questions et des variables sont faites soit dans les fichiers eux-mêmes soit dans des documents à part.

On a ainsi accès à des informations sur la position de la question dans le questionnaire, les filtres dans le questionnaire, les techniques de recodage, les consignes données aux enquêteurs, etc.

La documentation fine des questions et variables

https://explore.cdsp.sciences-po.fr/



La description papier de l'ensemble des variables d'un questionnaire

Etats-Unis

Emmanuel LAZEGA, Fondation Nationale des Sciences Politiques (FNSP)

Le phénomène collégial

Study Documentation

Table of Contents

Overview	
	wenge
Producers &	Sponsors
Sampling	
Data Collec	tion
	sing & Appraisal
	<u>Y</u>
	isclaimer
Files Descri	ption
	attributs.
	triplets.
	Multiplexite I
Variables L	
	attributs
	<u>triplets</u>
	Multiplexite I
Variables D	escription.
	attributs.
	triplets.
	Multiplexite_L

Files Description

Dataset contains 3 file(s)

tributs			
Cases	71		
Variable(s)	33		
ile Content	est le premier de cette série. Il représente, par avocat, l'ensemble des attributs et des variables strictement		

On retrouve ainsi dans ce fichier les propriétés socio-démographiques, les choix normatifs concernant l'avenir de la firme, ou encore les données comptables sur les performances économiques.

Ce fichier attribut peut être complété par le dossier compressé comportant les matrices conseil, amitié et collègues, nommé

Dans ce dossier, chacune des matrices, sous forme de fichier esv, présente en lignes et en colonnes l'ensemble des avocats. A l'intersection d'une ligne et d'une colonne est alors inscrit "0" ou "1" suivant qu'une relation existe ou non entre les deux avocats, si cette personne l'a déclaré dans son questionnaire.

En plus de ces trois matrices qui concernent l'ensemble du cabinet, il existe aussi une matrice "influence", qui correspond pour chacun des 36 associés aux relations d'écoute au sein de l'Assemblée générale des associés (qui écoute qui).

triplets	
# Cases	3043
# Variable(s)	46

Ces relations sont de type : "pour faire pression sur J. moi I, j'utiliserais K". Les informations présentes dans ce fichier permettent de mieux comprendre les relations au sein du cabinet.

La question posée était la suivante (Lazega et Lebeaux 1995) :

Voici la liste de tous les associés de votre cabinet. Imaginez que vous êtes le directeur (managing partner). Vous vous rendez compte que l'un de vos associés a des problèmes personnels qui ont des répercussions négatives sur sa productivité. Ces problèmes peuvent être de toutes sortes : alcoolisme, dépression, divorce, etc. En tant que directeur, c'est à vous de vous préoccuper de cette situation. Vous cherchez parmi les associés de cette personne en difficulté un ou des collègues qui iraient lui parler discrètement et confidentiellement pour savoir ce qui se passe, et pour voir ce que le cabinet peut faire pour aider et limiter les dégâts. Vous ne voulez pas le faire vous-même parce qu'il faut que la démarche reste informelle, et votre statut de directeur pourrait être gênant à cet égard. Ma question est la suivante : à qui, parmi tous les autres associés, demande riezvous d'aller parler à Associé No 1, si c'est lui qui est en difficulté ?

Les informations du jeu de données "triplets" sont aussi présentes au sein du dossier "pression_laterale". Ce dossier contient 36 matrices (une par associé) de 36*36 cases (levier * cibles). Chaque matrice représente les données brutes sur les choix de leviers (k) effectués par chaque répondant (i) pour faire pression sur chaque cible (j). L'information contenue dans ces matrices est aussi représentée dans le fichier « Triplet » où chaque observation équivaut à un choix de levier (n=3043). Dans la littérature ces bases de données s'appellent « three-way datasets ». Elles sont ici utilisées pour reconstituer et analyser le régime de contrôle latéral entre pairs. Pour cela, ces données sont enrichies par les attributs ou caractéristiques personnelles des acteurs et par la combinatoire des choix de leviers et de l'existence (ou non) de relation de collaboration, de conseil et/ou d'amitié entre répondant et levier, répondant et cible et levier et cible.

Apports et limites de la documentation "classique" des données

L'enquête en tant que processus de recherche global est documentée via des standards de métadonnées de type DDI, ELSST, etc.

Les descriptions détaillées des questions et des variables sont faites soit dans les fichiers eux-mêmes soit dans des documents à part ; on a ainsi accès à des informations sur la position de la question dans le questionnaire, les filtres dans le questionnaire, les techniques de recodage, les consignes aux enquêteurs, etc.

Beaucoup d'informations sont présentes.

L'inconvénient est qu'elles sont disséminées entre plusieurs documents / sources.

La globalité et des détails du processus de recherche doivent être recherchés dans les publications.

Apports et limites de l'"enquête sur l'enquête" (pour les entretiens)

Idée = rassembler, dans un document unique, les informations nécessaires à la compréhension & manipulation du corpus de données issues d'un processus de recherche

On est relativement proche du format data paper dans la structuration des rubriques & dans la forme littéraire

Cela donne des rapports de plusieurs dizaines de pages (2 à 3 fois la taille d'articles standards)

Page d'accueil de l'enquête sur l'enquête "Le choix de l'école"



Sommaire du rapport de l'enquête sur l'enquête "Le choix de

Enquête sur l'enquête « Choisir son école », beQuali, 2016,

l'école"

Sommaire	
INTRODUCTION	
1- GENESE DE L'ENQUETE	
1.1- Parcours de recherche	
1.2- QUEL POSITIONNEMENT FACE A L'OBJET DE RECHERCHE ?	
2- ANCRAGES THEORIQUES	1
2.1-LES REORIENTATIONS DE LA SOCIOLOGIE DE L'EDUCATION ET L'ETAT DES SAVOIRS SUI	R LE CHOIX
DE L'ECOLE	
2.2- LE CADRE THEORIQUE DU MODELE DES CHOIX SCOLAIRES	
2.3- CONSTRUCTION DE L'OBJET ET PROBLEMATIQUE DE RECHERCHE	10
3- REALISATION DES TERRAINS	19
3.1- L'ORGANISATION GENERALE DE LA RECHERCHE : UNE AGREGATION D'ENQUETES	19
3.2- * OBSERVER * LES CHOIX VERSUS RECUEILLIR LES DISCOURS DES ENQUETE(E)S	2
3.3- L'ORGANISATION DU TRAVAIL DE COLLECTE DES TEMOIGNAGES	2
4-CORPUS	2
4.1- LE CORPUS EXPOSE DANS CHOISIR SON ECOLE	2
4.2- LE CORPUS CONSERVE ET MIS A DISPOSITION	
4.3- RETOUR SUR L'ANONYMISATION	3
5-ANALYSE	3-
5.1-Retour sur la demarche d'analyse	3
5.2 LES ENTRETIENS ; QUEL CREDIT ACCORDER AU DISCOURS DES ENQUETE(E)S ?	3
5.3-LES PRINCIPALES INTERPRETATIONS PROPOSEES DANS CHOISIR SON ECOLE	3
6-POSTFACE	39
6.1-L'EXPLOITATION DE L'ENQUETE	
6.2. QUELLE GENERALISATION POSSIBLE DU MODELE DE CHOIX DE L'ECOLE ?	4
6.3-LES PISTES DE PROLONGEMENT ET DE REUTILISATION	4
BIBLIOGRAPHIE	40
OUVRAGES	4
Articles	4

Premiers bilans

On se situe dans le modèle du "faire pour - et - avec"

Un **coût d'entrée élevé** pour les ingénieurs extérieurs à l'enquête & un coût de documentation globalement élevé

Une **nécessaire collaboration** - même minimale - avec le chercheur déposant (éviter la mise à distance)

Les limites de l'exercice

- quelle posture adopter ? posture compréhensive versus posture critique
- jusqu'où aller dans la recontextualisation ?

Les freins de publicisation liés au format "littérature grise"

Partie 2

Retour sur la mise en place d'un service d'accompagnement à Sciences Po

La collection en auto-dépôt de data.sciencespo



Création d'un groupe de travail dédié

Un groupe multi-métiers (data librarians, documentalistes, ingénieurs BAP D)

But : mettre en capacité des ingénieurs de laboratoires, de la bibliothèque de Sciences Po, du CDSP

d'**orienter** et de **conseiller** les chercheurs (choix de revues cibles, aide ou participation à la rédaction de data papers)

<u>après</u> ou <u>au momen</u>t d'un dépôt d'enquête

Les linéaments et les contours de ce projet sont décrits dans un chapitre à paraître dans les actes du colloque #DHNord 2021 sur les data papers

Les premières actions menées

Etat des lieux des pratiques en vigueur dans l'écosystème éditorial impliquant la science politique et la sociologie principalement

A partir d'un **échantillon raisonné** de data journals anglophones ou des revues académiques plus classiques (disciplinaires ou méthodologiques) en France ou dans le monde anglophone

Démarche incrémentale : principales revues dans lesquelles les chercheurs de Sciences Po publient & des Data journals ou revues disciplinaires qui se font fait connaître sur le sujet

Principe : aider à **déterminer la cible** aisément (positionnement de la revue, coûts, particularités du processus d'évaluation, etc.)

Articulation de 2 logiques :

- une logique réactive : conseiller les (rares) chercheurs déjà informés et motivés
- une **logique proactive** : se mettre en capacité de promouvoir cette possibilité auprès de chercheurs plus distants

Les actions en cours

Mise en place de tutoriels / guidelines

- "scripts" sur les plans, informations demandées, points d'attention, check lists... permettant d'alerter sur les besoins d'élaboration & de rédaction
- mise en place de "modèles" de DP servant de cadres de référence
- sert de "maïeutique" dans un échange entre chercheur et ingénieur

Mise en place de tests en grandeur nature avec des "cobayes"

Les limites principales de cette démarche : **encore peu de débouchés** (gratuits) dans les revues disciplinaires ou les Data Journals, même si dans le même temps il y a encore peu de candidats

⇒ quel degré d'adéquation entre l'offre et la demande ?

Conclusion

Quels enjeux autour des data papers?

Un flou sur les contours éditoriaux des data papers

Une difficulté à **formaliser un modèle de data paper** pour les disciplines visées, la différence et l'articulation avec la forme "article de recherche" n'étant pas encore très lisibles

Une question non résolue : le **statut de cette publication**, de l'évaluation dont elle peut faire l'objet, en lien avec le ou les articles scientifiques consécutifs de la même recherche

Question du **support de publication**, des rôles respectifs des revues classiques et des data journals, généralistes ou spécialisés, avec un risque d'effet de niche

Une question souvent éludée : quel(s) lectorat(s) cible(s) des data papers, en lien avec les pistes de réutilisation fournies par les auteurs - recherche versus pédagogique

Un flou sur le périmètre méthodologique des data papers

Un **tropisme "quanti"**, lié aux différences d'ancrage de la pratique de partage / réutilisation dans les différentes cultures professionnelles

Quelle place pour les données "quali" ?

Un point aveugle : quid des données du Web ?

Une question centrale mais souvent évitée : les coûts humains de la rédaction des DP

Quel **bilan coûts / avantages** impliqués par la production des data papers, pour les chercheurs ou les ingénieurs en soutien ?

Quelles formes de division du travail entre chercheurs et personnels IST (archivistes, data librarians, documentalistes, ingénieurs Bap D)?

Vers un changement du **positionnement des ingénieurs** face aux chercheurs en matière d'ouverture des données de la recherche ? en particulier s'agissant de la nature des services qu'ils rendent

Merci de votre attention

<u>anna.egea@sciencespo.fr</u> <u>guillaume.garcia@sciencespo.fr</u>