

Introduction to effective and practical research data management

Cristina Magder

Data Collections Development Manager

06 December 2024



Format of the workshop

- Presentation
 - Explore key concepts and principles of research data management and gain insights into best practices and strategies for effective data handling with plenty of tools and additional resources available.
- [Mentimeter](#)
 - Engage with practical exercises to reinforce concepts. Share your thoughts, ideas, and responses in real-time.
- Q&A
 - Conclude with a live Q&A session during the workshop, please use the Zoom Q&A function to ask any and all questions you have about data management.

Today's programme

- Research data lifecycle.
- FAIR and CARE data principles.
- Data management planning.
- Overview of ethical and legal considerations.
- Data security, storage and backup.
- Data curation best practices.
- Data sharing strategies.
- Q&A session.



Learning objectives

Gain an overview of the research data lifecycle, data management planning, and the principles underpinning good data practices, including FAIR and CARE guidelines.

Learn about effective methods to organise, document, store, and share research data responsibly, ensuring compliance with ethical and legal requirements.

Effective research data management

Effective research data management (RDM) practices ensures data are:

- compliant with ethical standards and applicable legislation
- well-organised, quality controlled and well-documented
- safely stored, backed up, processed and analysed
- responsibly archived and preserved, and appropriately shared for future reuse

RDM practices safeguard the integrity of the research.

“UK Research and Innovation (UKRI) expects research data arising from its funding to be made as open as possible and as restricted as necessary. Good research data management practices should be followed throughout your project.”

[Source](#)

Case Study: Repurposing inspection data for research – the HMIP survey journey

- His Majesty's Inspectorate of Prisons (HMIP) collects data during inspections to monitor prison conditions and treatment.
- By repurposing this data for research, it has been possible to provide broader insights into prison reform and policy.
- Data was effectively documented, archived, and shared for secondary use.
- Researchers used the data to explore themes like institutional environments and prisoner wellbeing.
- Comprehensive metadata and proper anonymisation ensured the data's usability while respecting confidentiality.

[Read more about this.](#)

Case Study: The data life cycle of an archived qualitative study used for teaching

- A qualitative study, originally conducted to explore social behaviours, was preserved in a data archive and later repurposed as a teaching resource.
- Detailed documentation and careful curation ensured the study's relevance long after its original use.
- The archived data has been used to teach students about research methodologies, data analysis, and ethical considerations.
- Effectively managed data can serve multiple purposes, maximising its educational and research value.

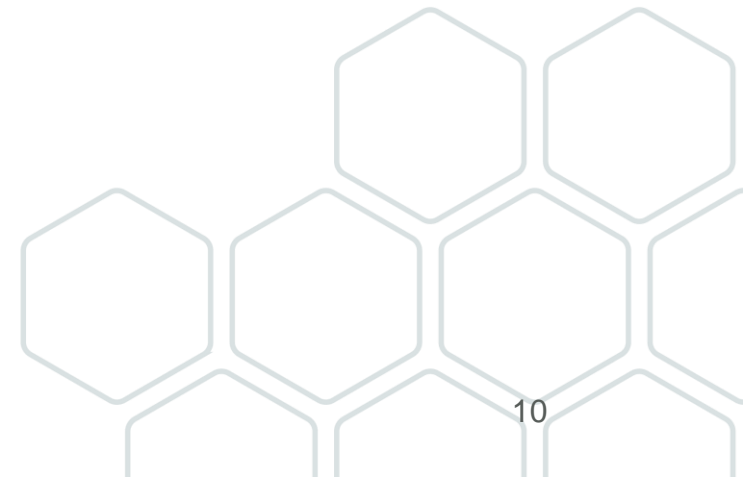
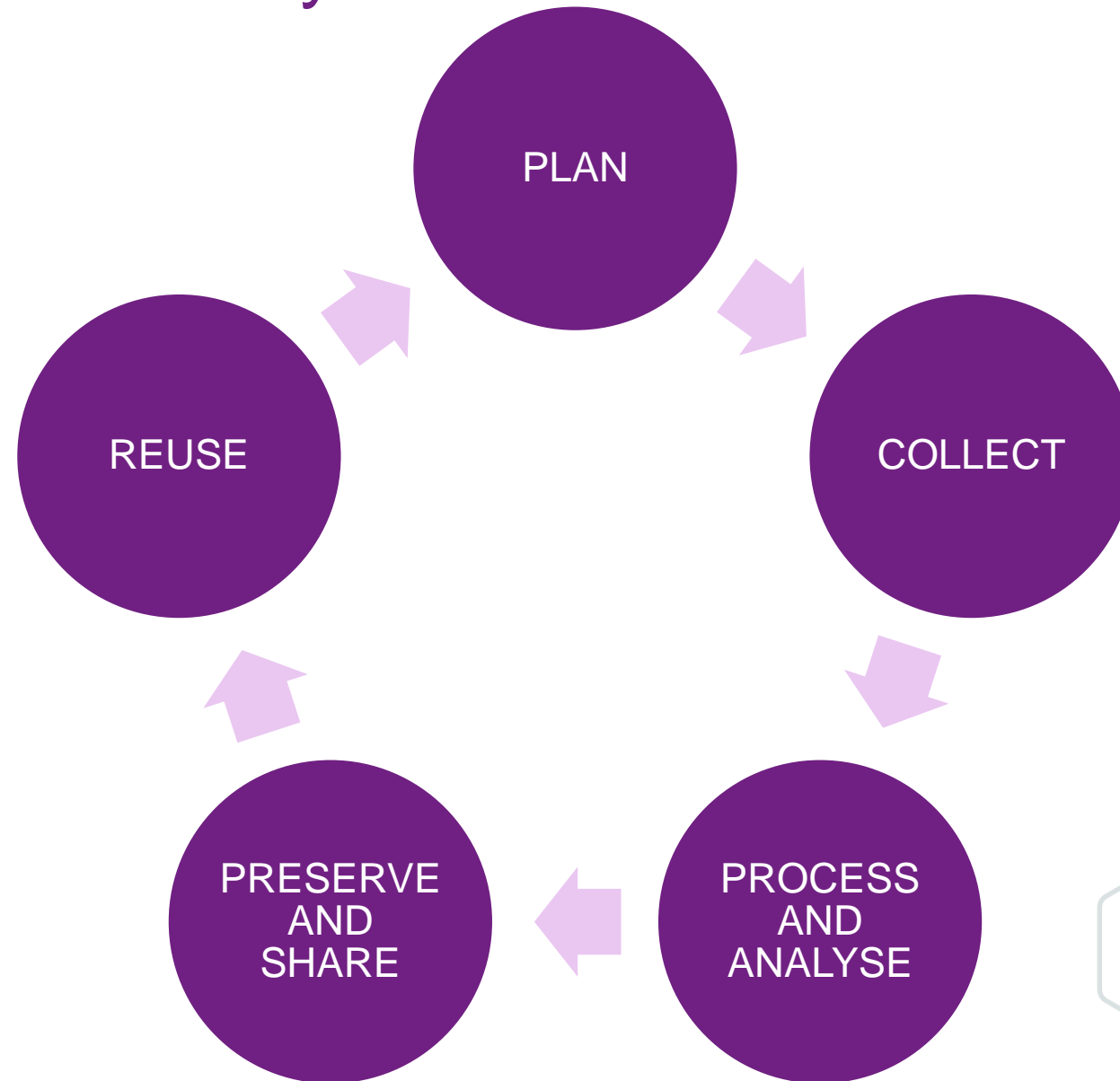
[Read more about this.](#)

Case Study: Balancing teens' privacy with the desire to share data

- A study involving data about teenagers had to address ethical and legal challenges to ensure participants' privacy was protected while allowing data sharing.
- Sensitive information was anonymised to protect individual identities.
- Researchers balanced privacy with the need for transparency and data reuse.
- Adhering to ethical guidelines and employing anonymisation techniques enabled appropriate data sharing without compromising participant trust.

[Read more about this.](#)

Research data lifecycle



FAIR data principles

Findable

Accessible

Interoperable

Reusable

Published in 2016 in *Scientific Data*, FAIR principles are outlined as guidelines to help define what good data management looks like, emphasising machine-actionability.

Find out more about the [FAIR initiative](#).

How FAIR aware are you?

Explore the DANS FAIR-Aware assessment tool to evaluate and improve your understanding of the FAIR principles for managing research data effectively.

Answer a series of questions to reflect on your current practices.

Receive feedback to help align your data management with FAIR principles.

[FAIR-Aware Assessment Tool](#)

How to be FAIR

Findable

- community-endorsed discovery metadata standards
- machine readable discover metadata
- unique persistent identifiers e.g. DOIs.

Accessible

- data licensing and availability statements (including restrictions)
- methods/tools to access the data
- metadata preserved indefinitely.

Interoperable

- standard vocabularies/ontologies
- standard metadata schemas.

Reuseable

- community-endorsed data licensing
- provenance information in metadata
- established data quality assurance processes
- open format files for long-term preservation.

#BeFAIRandCARE



[Read more on CARE \(the Global Indigenous Data Alliance\)](#)

Even more success stories

Using [Understanding Society](#), researchers looked at the way engaging with arts, culture and sports can [lead to greater satisfaction with life](#).

Studies like the [English Longitudinal Study of Ageing](#) can give insight into how [enjoying later life can be linked to living longer](#) and the impact of [social isolation and loneliness on mortality in older people](#).

The [Millennium Cohort Study](#) was used to [investigate the conditions associated with parental involvement with children](#) for policy recommendations.

[Family Resources Survey](#) to [analyse migrants' experiences of poverty and to compare them with the experiences of UK-born people](#).

Data Management Plans (DMPs) overview

DMP topics overview

From a generalised perspective DMPs should cover:

1. Data description: new and existing data.
2. Ethical and legal considerations and compliance.
3. Curation of data: organising, formatting, and documenting.
4. Data security, storage and backup.
5. Data sharing strategies.
6. Responsibilities and resources.

We provide in-depth guidance on our website for the [ESRC DMP](#).

Always remember to check funder requirements and that a DMP is a living document, as research evolves the DMP should be reviewed and updated as necessary.

Why is data management planning essential?

- Anticipate and prepare.
- Keep on track.
- Secure necessary resources.
- Think ahead about storage and safeguard data.
- Share FAIR data and ensure reproducibility.
- Meet funders' expectations.

DMPonline

DMPonline - web-based platform developed by the Digital Curation Centre to help researchers create, review, and share DMPs that meet institutional and funder requirements

The screenshot shows the DMPonline web interface. At the top is an orange navigation bar with the DMPonline logo, 'My Dashboard', 'Create plans', 'Reference', and 'Help'. A 'Language' dropdown is on the right. Below the navigation bar is a grey header with 'University of Essex'. A light blue info box contains a message about SSO login upgrades and instructions for re-linking accounts. Below this is the 'Create a new plan' section, which includes a form for entering research project details, selecting a primary research organisation (University of Essex is selected), and selecting a primary funding organisation. The form has 'Create plan' and 'Cancel' buttons at the bottom.

Info:
As part of our routine maintenance, we have upgraded our SSO login to enhance security.

- If your account was not linked to your institutional credentials, please log in as normal.
- If your account was linked to your institutional credentials, you will now need to re-link your account as part of this upgrade. To do this, please log in using your **DMPonline** email and password.
- Next, go to **Edit profile** > scroll down to the point **Institutional credentials**, and select the **Link your institutional credentials** option.
- After re-linking your account, you need to refresh your browser to complete the process. Remember to save your updated settings.

Create a new plan

Before you get started, we need some information about your research project to set you up with the best DMP template for your needs.

* What research project are you planning?

mock project for testing, practice, or educational purposes

* Select the primary research organisation

Organisation - or - No research organisation associated with this plan or my research organisation is not listed

* Select the primary funding organisation

Funder - or - No funder associated with this plan or my funder is not listed

Ethical and legal considerations overview

Ethical considerations

- Maximise benefits and minimize risks.
- Voluntary participation and informed consent.
- Respect individual rights and dignity.
- Integrity and transparency.
- Clear responsibilities and independence.

Check our [Role of Informed Consent in Ethical Data Collection, Sharing and Reuse workshop materials](#)
(presentation and recording available)

Legal considerations

- Lawful data handling.
- Data minimisation and anonymisation.
- Secure storage and access controls.
- Data retention and disposal.
- Data sharing protocols.
- Fairness, transparency and accountability.
- Intellectual property rights.

Check our [Ethical and legal guidelines in data sharing workshop](#) materials (presentation and recording available)

How do ethical and legal considerations in data management affect the integrity and impact of research?



- Protect participants.
- Maintain and build trust.
- Enable sharing and collaboration.
- Ensure long-term impact.
- Enable appropriate risk management.

Data security, storage and backup

Data security and storage



Data must be protected data from unauthorised:

- Access.
- Use.
- Change.
- Disclosure.
- Destruction.

Digital back-up strategy

Control access to computers:

- Use pasphrases and lock your machine when away from it.
- Run up-to-date anti-virus and firewall protection.
- Power surge protection.
- Restrict access to sensitive materials e.g. consent forms.
- Always keep personal data separate, secure and encrypted.
- Utilise encryption:
 - on all devices: desktops, laptops, memory sticks and mobile devices.
 - at all locations: work, home and travel.

Control physical access to buildings, rooms and filing cabinets.

Properly dispose of data and equipment.

Digital back-up strategy

- Making backups of files is an essential element of research data management which ensures that original data files can be restored from backup copies, should they get damaged or go missing.

Three+ copies of the data, with at least one being stored offsite.

- Regular backups help protect against **accidental** or **malicious** data loss due to:
 - human error
 - hardware failure
 - software or media faults
 - virus infection or malicious hacking
 - power failure.

[Further information on storing data.](#)

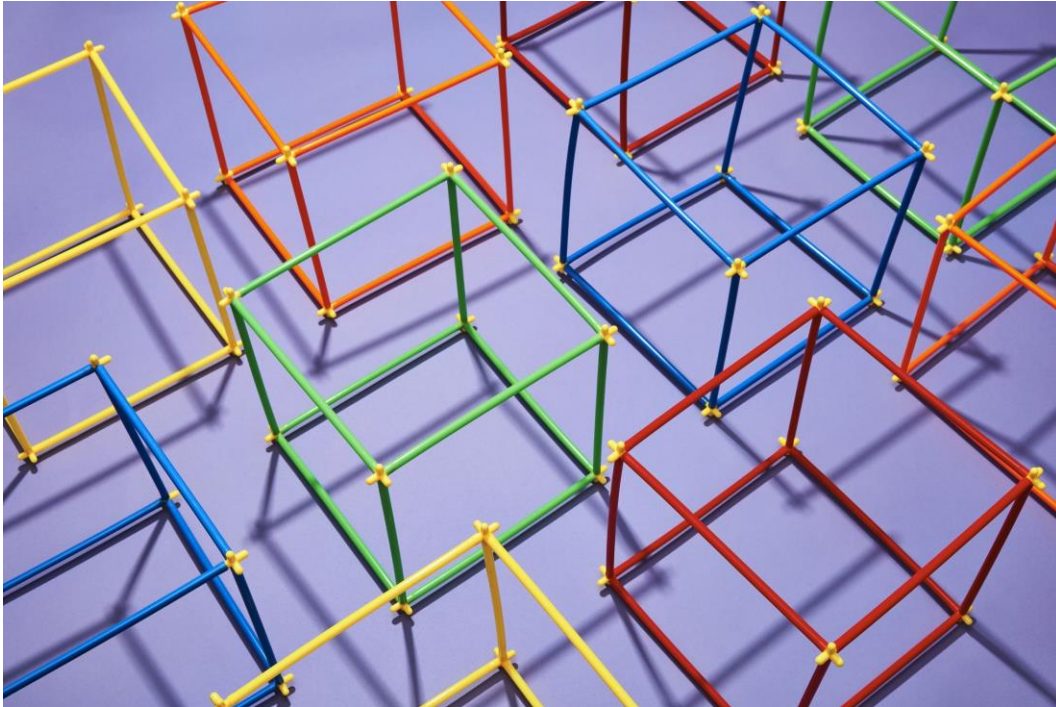
Data disposal

- Simply deleting files and reformatting a hard drive will not securely erase information, meaning that it will still be possible to recover the data that was previously on the hard drive.
- Software is available to help erase files from hard disks, meeting recognised erasure standards. Example software is: [BCWipe](#), [WipeFile](#), [DeleteOnClick](#) and [Eraser](#) for Windows platforms; and [Permanent Eraser](#) for MacOS platforms.
- Shredders certified to an appropriate security level should be used for destroying paper and optical media.

[Further information on disposing data.](#)

Data curation: formatting organising and anonymising

File formats strategy



[Further information on formatting data.](#)

What format is best suited for data creation?

What format is best suited for data analyses and other planned uses?

What format is best suited for long-term sustainability and sharing of data?

Should you choose an open versus a proprietary format?

Should the format be lossy or not?

Is the format suitable for conversion?

Recommended file formats

Type of data

Quantitative tabular data with extensive metadata.

A dataset with variable labels, code labels, and defined missing values, in addition to the matrix of data.

Qualitative data. Textual.

Recommended formats

Proprietary formats of statistical packages e.g. SPSS (.sav), Stata (.dta), .sas7bdat.

Delimited text and command ('setup') file (SPSS, Stata, SAS, etc.) containing metadata information.

Some structured text or mark-up file containing metadata information, e.g. DDI XML file.

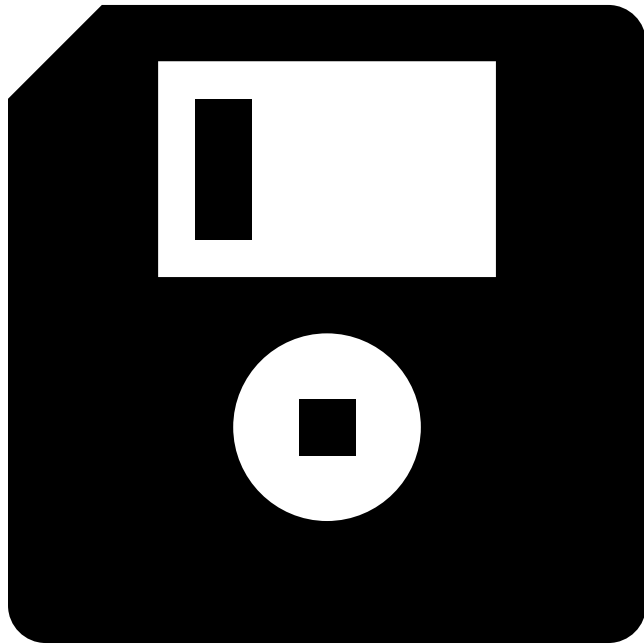
eXtensible Mark-up Language (XML) text according to an appropriate Document Type Definition (DTD) or schema (.xml).

Rich Text Format (.rtf).

Plain text data, ASCII (.txt).

[UK Data Service recommended formats](#)

Best practices for version control



Data management processes inevitably create a number of edits to the data and documentation.

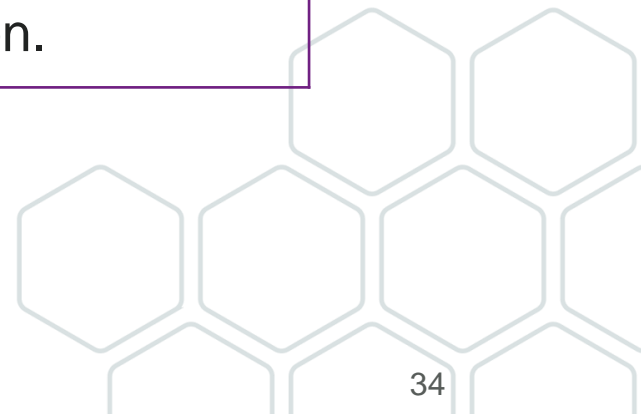
- Identify milestone versions to keep
- Uniquely identify different versions
- Record changes
- Record relationships between items
- Track the location of files
- Regularly synchronise files
- Identify a single location for the storage of milestone and master versions.

Best practices for file naming

- Develop a file naming strategy (minimal elements version number, description of content, publication date, project number).
- Create meaningful but brief names.
- Use file names to classify types of files.
- Use dates in the format YYYY-MM-DD.
- Avoid capitalisation where possible, as some computer platforms may be case-sensitive (e.g. Unix)
- Avoid using spaces, dots and special characters (& or ? or !).
- Use hyphens (-) or underscores (_) to separate elements in a file name.
- Reserve the 3-letter file extension for application-specific codes of file format (e.g. .docx, .xlsx, .mov, .tif).
- Include versioning within file names where appropriate.
- Review file names for archived versions to ensure they do not contain any confusing, irrelevant information (e.g. versioning, misleading description).

File naming examples

File name	Meaning
pn6614_ukhls_wave2_e10_2024-04-15.sav	Project number 6614, Wave 2 of the UKHLS SPSS data 10 th edition, last edited on 15 April 2024.
pn018_int127_js_v1_2024-03-02.rtf	Project number 18, transcript of interview with participant 127, conducted by JS on 2 March 2024, first version.
shes_21_dataset_documentation_v2.pdf	Scottish Health Survey 2021 dataset documentation second version.



Variable formatting and measurement levels

Incorrect variable formatting can lead to incorrect data use and sometimes even reduces the usability of the data.

A very simple but important is to **determine whether the data are to be treated as string or numeric.**

Numeric variables also need to be checked to ensure the measurement level is correctly defined. Note that different software name and treat measurement levels differently.

Variable measurements examples

Variable	Stata	SPSS
Ethnicity	Categorical	Nominal
Annual income (banded)	Categorical	Ordinal
Marital status	Categorical	Nominal
Age (banded)	Categorical	Ordinal
Monthly income (£)	Continuous	Scale

STATA: Categorical, Continuous
SPSS: Nominal, Ordinal, Scale



Quality assurance

You must **check and document** any changes made to your data. This will provide a history, version control and provenance trail to help quality assure your data

Quality assurance checks may include:

- Double-checking coding of observations or responses and out-of-range values.
- Checking data completeness.
- Adding variable and value labels where appropriate.
- Statistical analyses, such as frequencies, means, ranges or clustering to detect errors and anomalous values.
- For qualitative interview data, correcting errors made during transcription.

QAMyData

[An open-source tool](#) developed by the UK Data Service to perform 'health checks' on numeric data. It automates the detection of common issues such as missing values, duplicates, outliers, and direct identifiers, ensuring your datasets are accurate and reliable.

- Identifies missing data, duplicates, outliers, and potential direct identifiers.
- Customise assessments to align with your project's specific data quality standards.
- Generates detailed summaries highlighting areas that may require attention.

Tools for qualitative and linguistic data

CLARIN-NL: A hub for linguistic data and tools integrated into the European CLARIN infrastructure; linguistic datasets, analysis tools, demonstrators, and applications.

[Explore CLARIN-NL.](#)

CLARIN Switchboard: An intuitive platform linking datasets to suitable CLARIN tools; suggests tools for tasks like tokenisation, parsing, and annotation.

[Explore CLARIN Switchboard.](#)

ELAN: A flexible tool for annotating and analysing audio and video data; multi-tiered, time-aligned annotations, ideal for working with interviews, focus groups, or observational data

[Explore ELAN.](#)

Three-prong approach to protecting participants

- **Consent:** participants must be informed about risks and benefits of *any* data sharing.
- **Anonymisation:** treat the data reducing the risk of re-identification recognising data utility will be reduced, therefore a balanced approach must be taken
- **Access:**
 - Who? How? For how long?
 - Access levels and user agreements
 - Leverage legislation and existing frameworks such as the Five Safes Framework

Check our [Introduction to anonymisation techniques for social sciences research data workshop](#) materials (presentation and recording available)

Useful semi-automated anonymisation tools

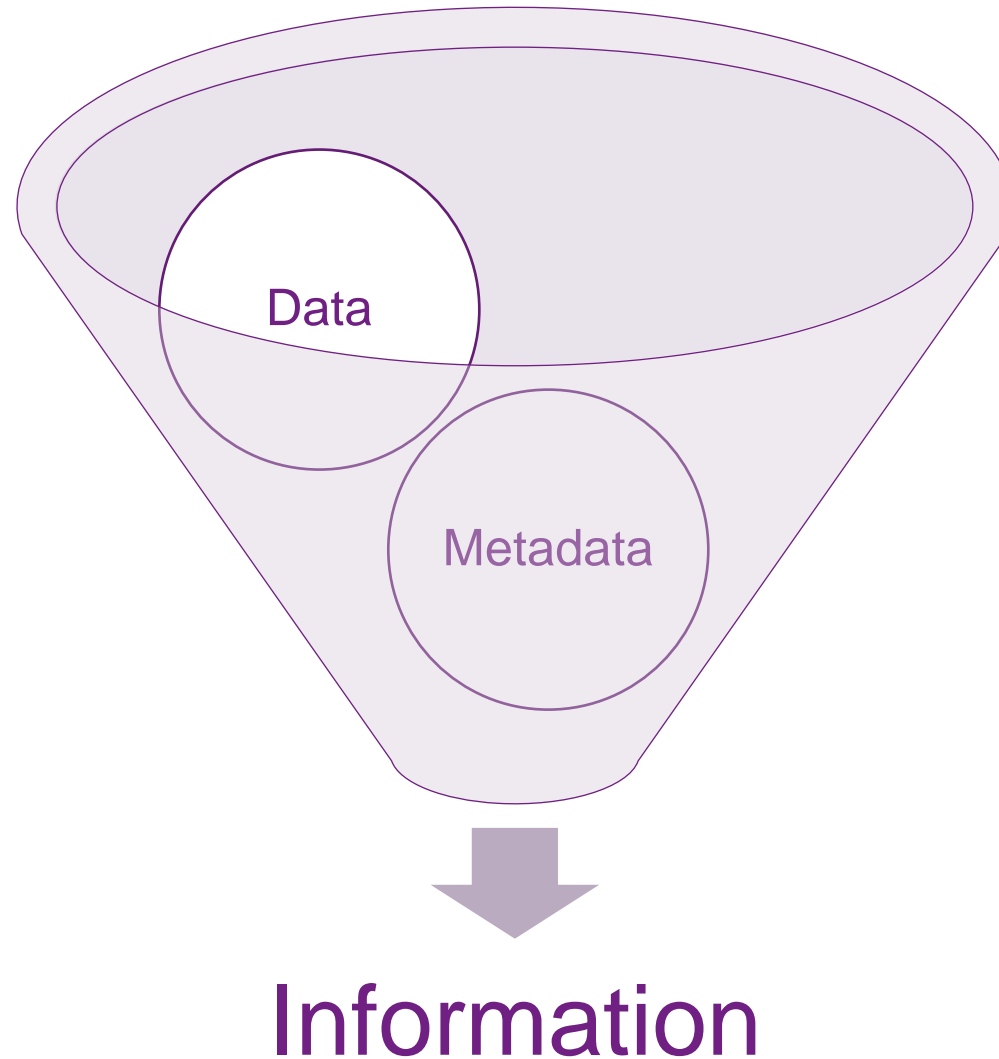
- [sdcMicro](#) – R package (free) – has a user-friendly interface so minimal coding skills needed.
- [QAMyData](#) - UK Data Service developed a free (GitHub) easy-to-use open source tool, that provides a health check for numeric data. The tool uses automated methods to detect and report on some of the most common problems in survey or numeric data, such as missingness, duplication, outliers and direct identifiers.
- [ARX](#) - a comprehensive open-source software for anonymizing sensitive personal data. It supports a wide variety of (1) privacy and risk models, (2) methods for transforming data and (3) methods for analysing the usefulness of output data.
- [μ-Argus](#) – developed by Statistics Netherlands; [User Manual](#)
- [QuailAnon](#) – open-source tool developed by QualiService
- [Text anonymisation helper tool](#) – Word Macros tool developed by the UK Data Service
- [Textwash](#) - open-source tool uses Python to identify and replace direct identifiers
- [FAMTAFOS](#) – in development; open-source desktop app that utilises AI technology to anonymise text at scale; operates on principle of Named Entity Recognition (NER), and can be set to search for names, locations, occupations, etc. They will then tag them for subsequent human editing.
- [De-ID](#) - HIPPA-compliant tool to flag potentially identifiable data; only available to organisations

Data curation: metadata and documentation

Metadata

For data to become information, you need to understand the context in which the data are situated.

Metadata is what provides this essential context.



Why is it essential to document data?

- Efficiency and accessibility.
- Ethical and accurate data reuse.
- Reproducibility and validation.
- Compliance and ethical standards.
- Long-term preservation.

Check our [Best practices for documenting social sciences research data](#) (presentation and recording available)

Types of data documentation

Data-level documentation

- provides information on the individual data objects, such as a variable in a data file or an interview transcript. It can be embedded in the data file, such as variable or value labels in a data file, or participant information added in the header or an interview transcript.

Study-level documentation

- provides high-level information on the research context and design, the data collection methods used, any data preparations and manipulations, plus summaries of findings based on the data.

[Further information on documenting data.](#)

Example for survey data

	Name	Type	Width	Decimals	Label	Values
1	serial_scr...	Numeric	8	2	Scrambled Individual Serial	None
2	BSA21_fina...	Numeric	8	2	British Social Attitudes Survey 2021 - final weight	{-9.00, Refu...
3	DVSex21	Numeric	3	0	What is your sex?	{-1, Not app...
4	Ragecat	Numeric	2	0	Age of respondent(grouped) <7 category> dv	{-1, Not app...
5	HHincomex	Numeric	8	2	What is the total income of your household from all sources before tax?	{-1.00, Not ...
6	RClassGP	Numeric	1	0	NS-SEC analytic classes (self-coded)	{-1, Not app...
7	hedqual2x	Numeric	8	2	Highest educational qualification attained	{1.00, Degr...
8	MarStat6x	Numeric	8	2	Marital Status	{1.00, Marri...
9	HhICHlGpdx	Numeric	8	2	Children in household (grouped)	{1.00, Yes}...

Value Labels

Spelling...

Value Labels:

Value	Label
1.00	Married/in a civil partnership/living with a partner
2.00	Separated/divorced/dissolved civil partnership
3.00	Widowed/surviving partner from a civil partnership
4.00	Single (never married/never in a civil partnership)

OK Reset Cancel Help

University of Manchester, Cathie Marsh
Institute for Social Research (CMIST), UK
Data Service. (2024). British Social Attitudes
Survey, 2021, Health Care and Equalities:
Open Access Teaching Dataset. [data
collection]. NatCen Social Research,
[original data producer(s)]. NatCen Social
Research. SN: 9236, DOI:
<http://doi.org/10.5255/UKDA-SN-9236-1>

Example for interview data

Name / Interview ID	Year of Birth	Place of Birth	Gender	Interviewer	Place of Interview	Number of occasions interviewed	Date of Interview	Full interview/ Summary	No of Pages	Text File Name	Qualbank link to interview	YouTube Playlist of audio clips
Frank Bechhofer	1935	Germany	Male	Paul Thompson	Edinburgh	1	8 January 2001	Interview	78	6226int001	https://discover.ukdataservice.ac.uk/QualiBank/Document/?id=q-a4e16904-7299-49f9-9b0a-4b601384a665	https://youtu.be/UaKqnUC4IbA
								Highlights	24	6226themext001	https://discover.ukdataservice.ac.uk/assets/qualibank/6226themext001.pdf	
								Summary	8	6226intsum001	https://discover.ukdataservice.ac.uk/QualiBank/Document/?id=q-d011d67f-63a0-436f-8e35-67cf677c59ca	
Daniel Bertaux	1941	France	Male	Paul Thompson	France	1	20 August 2002	Interview	82	6226int002	https://discover.ukdataservice.ac.uk/QualiBank/Document/?id=q-d9af5999609749c59fa01b594f714397	https://youtu.be/dxC4fz2zqzA
								Highlights	22	6226themext002	https://discover.ukdataservice.ac.uk/assets/qualibank/6226themext002.pdf	
								Summary	14	6226intsum002	https://discover.ukdataservice.ac.uk/QualiBank/Document/?id=q-2b777f6c915c422f867fe1c13683276b	
Mildred Blaxter	1925	Newcastle-on-Tyne	Female	Paul Thompson	Norwich	1	2 August 2002	Interview	57	6226int003	https://discover.ukdataservice.ac.uk/QualiBank/Document/?id=q-1a09ab2d-e529-4629-9201-4a103644111c	https://www.youtube.com/watch?v=R6G5e15sG68&list=PLaz84-Ixz3kiK5bOD5KBVAwe2eks_2GP
								Highlights	10	6226themext003	https://discover.ukdataservice.ac.uk/assets/qualibank/6226themext003.pdf	
								Summary	6	6226intsum003	https://discover.ukdataservice.ac.uk/QualiBank/Document/?id=q-5d116f52-c369-4ecc-96ad-efca438a2077	
Pat Caplan	1942	Neston, Cheshire	Female	Paul Thompson	London	1	18 November 2009	Interview	66	6226int004	https://discover.ukdataservice.ac.uk/QualiBank/Document/?id=q-51c8f358-2ce1-4012-b00f-87ce4bbe5a50	https://www.youtube.com/watch?v=HUAfRPOwxEw&list=PLaz84-Ixz3kiK5bOD5KBVAwe2eks_2GP
								Highlights	23	6226themext004	https://discover.ukdataservice.ac.uk/assets/qualibank/6226themext004.pdf	
								Summary	9	6226intsum004	https://discover.ukdataservice.ac.uk/QualiBank/Document/?id=q-84e8c183-29ce-481a-a5c6-b78f7ac92723	

Thompson, P. (2019). *Pioneers of Social Research, 1996-2018*. [data collection]. 4th Edition. UK Data Service. SN: 6226, DOI: <http://doi.org/10.5255/UKDA-SN-6226-6>

Study-level documentation example

Scottish Crime and Justice Survey, 2019-2020

Details Documentation Resources Access data

Documentation

Title	File name	Size (MB)
Scottish Crime and Justice Survey 2019-20: Disclosure Control Report	8799_scjs_ukda_disclosure_control_report_2019-20.pdf	0.45
Scottish Crime and Justice Survey 2019-20: Questionnaire and User Notes	8799_scjs_questionnaire_and_user_notes_2019-20.pdf	2.5
Scottish Crime and Justice Survey 2019-20: Technical Report	8799_scjs_technical_report_2019-20.pdf	2.25
UK Data Archive Citation File for Study 8799	UKDA_Study_8799_Information.htm	0
UK Data Archive Data Dictionaries	ukda_data_dictionaries.zip	0.12
UK Data Archive ReadMe File for Study 8799	read8799.htm	0

ScotCen Social Research. (2022). Scottish Crime and Justice Survey, 2019-2020. [data collection]. UK Data Service. SN: 8799, DOI: <http://doi.org/10.5255/UKDA-SN-8799-1>

Data sharing strategies

Data sharing strategies overview

There are various way of sharing research data including

- Domain specific repositories (data service providers, archives, centres)
- Institutional repositories
- Self-preservation and dissemination
- Commercial data sharing platforms
- Direct submission with journal publications

All data sharing methods have advantages and disadvantages.

Key consideration for data sharing strategies

Purposes of data sharing

Data findability, accessibility, interoperability and reusability

Data sensitivity and confidentiality

Ethical and legal implications

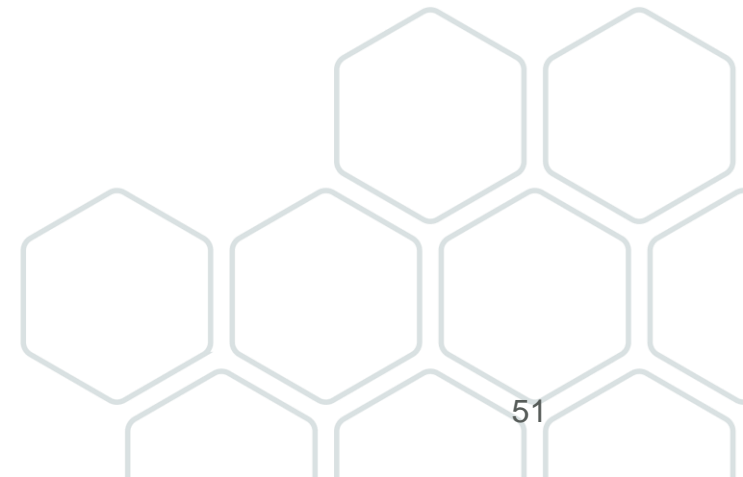
Long-term preservation

Security measures

Costs and resources

Technology and infrastructure

Stakeholder engagement and support



Responsible repositories

Responsible repositories are data service providers that host facilities that adhere to established standards and best practices in data management.

These ensure the integrity, preservation, and accessibility of the data they hold.

Responsible repositories play a critical role in the research ecosystem by providing reliable, secure, and accessible environments for storing and disseminating research data.



Source: [Research Data Alliance](#)
[TRUST principles](#)

Deposit licence agreements



Protect the rights of the data owner and the repository.



Ensure that data users are aware of their rights and responsibilities.



Facilitate ethical and legal sharing and use of data, enhancing its value to the research community.

What if I use secondary data?

Always check the licence under which the data are made available.

While you might not be able to share derived data you can always share your code.

 code/syntax file are clean, well formatted and do not contain any data

 avoid including unnecessary personal information

 code/syntax file are well commented

 always include the full citation (including the persistent identifier) for the data used

 provide in-depth metadata describing the files and methods used

 provide a ReadMe/Methods document for ease of use for secondary users

Registry of data repositories

Re3data.org is a comprehensive registry of research data repositories that has been active for over a decade. It provides a curated index of more than 3,000 repositories worldwide, covering all academic disciplines.



Source: [re3data](https://re3data.org)

Other tools, templates and tutorials

- [Research data management learning hub](#)
- [Data management checklist](#)
- [Data management costing tool and checklist](#)
- [Model consent form and survey consent statement](#)
- [Transcription template](#)
- [Transcription instructions](#)
- [Data list template](#)
- [Data skills modules](#)
- [UKDS events](#)

Get connected

[UK Data Service](#)

[Jisc mail group](#)

[@UKDataService](#)

[UK Data Service YouTube channel](#)

Powerpoint slides will be available on our website in due course and you can catch up on the recording on our Youtube channel.

Thank you ever so
much!

`datasharing@ukdataservice.ac.uk`

<https://beta.ukdataservice.ac.uk/help>

Any questions?

