# Data Handling:

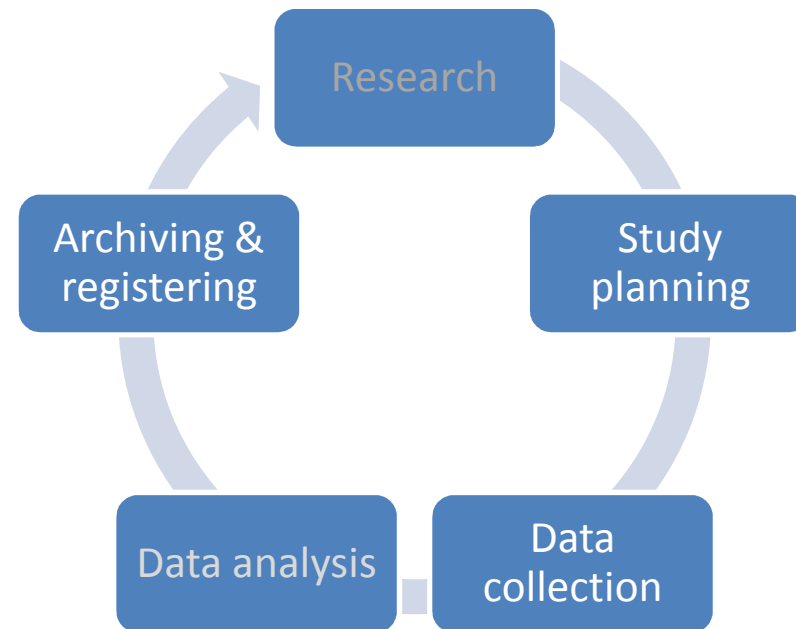# Documentation, Organization and Storage

Sebastian Netscher

CESSDA Training at the Data Archive for the Social Sciences

GESIS - Leibniz Institute for the Social Sciences

@CESSDA_Data

# Data Documentation

# Why Data Documentation?

# Levels of Data Documentation

- Study level
  - study description
  - study design
  - data processing

- Variable level
  - questionnaire
  - variables and codes



Image by A. Herrema & H. Bouwteam (CC-by)

# Structured versus Unstructured Metadata

**Unstructured documentation**

- technical reports etc.
- questionnaire, show cards, interviewer instructions etc.
- codebook etc.

**Standardized forms for standardized information**

- coding schemas, e.g. ISCO, ISCED etc.
- international metadata standards, e.g. DDI

# The Data Documentation Initiative (DDI)

- International standard for the description of data
  - DDI-Codebook (DDI2)
    ⇒ based on the codebook
  - DDI-Lifecycle (DDI3)
    ⇒ based on the (DDI) data lifecycle



Source: http://www.ddialliance.org/

# Persistent Identifiers (PIDs)

- Persistent identifiers
  - provide permanency
  - assure unique retrieval of data
  - assign citation for reuse

- The DOI system
  - controlled by IDF
    (*International DOI Foundation*)
  - DOI Resolver,
    e.g. http://www.doi.org/index.html

DOI: 10.ORGANISATION/ID
Example: 10.4232/1.11159
prefix/suffix

# da|ra Schema: Main Categories

Publication Date

Availability (controlled)

Notes

Sampled Universe

Relation

Contributor

Title          Other Titles

Dataset

Title

Creator          Description

Availability (free)

Version

Temporal Coverage

Alternative Identifier

Language

Classification Internal

Collection Mode (controlled)

Geographic Coverage

Rights                    URL
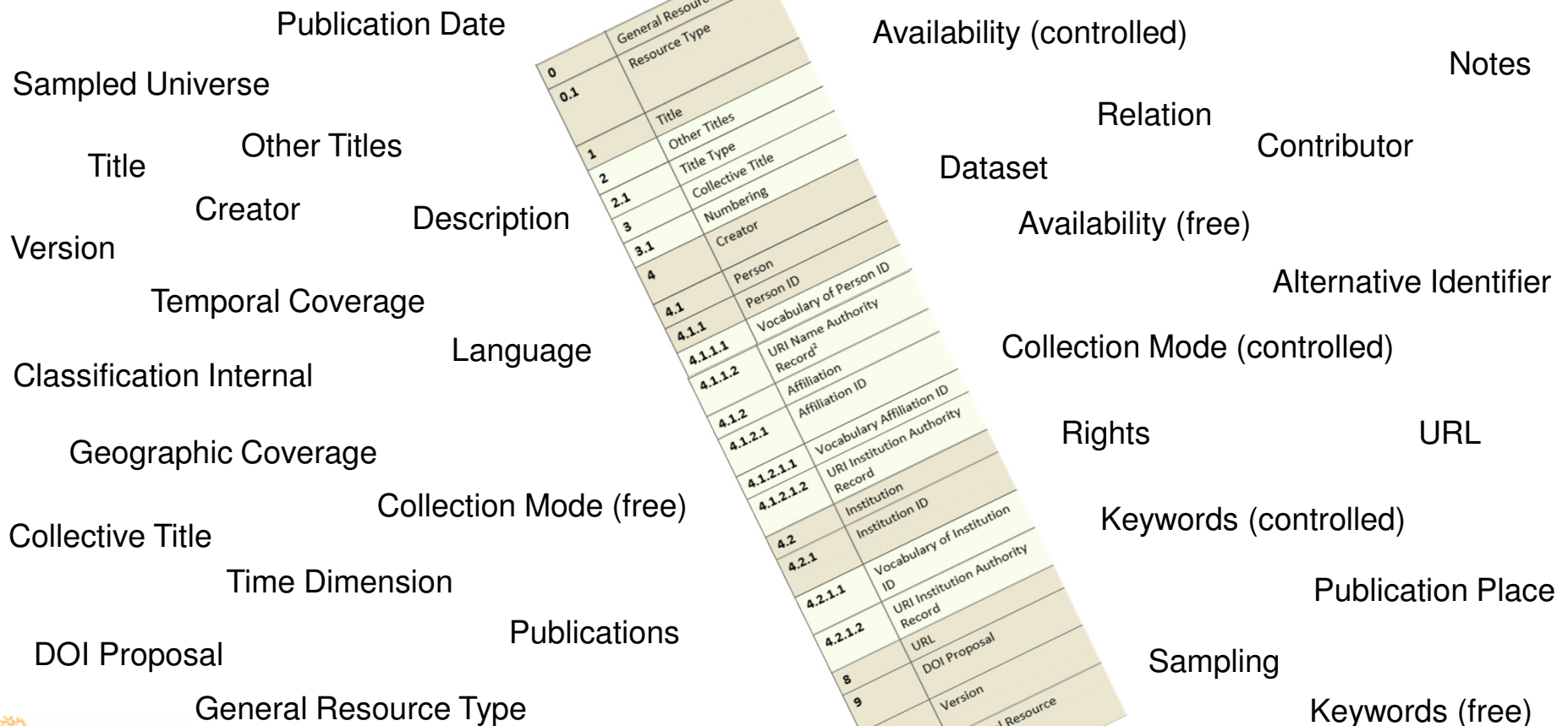
Collection Mode (free)

Collective Title

Keywords (controlled)

Time Dimension

Publication Place

Publications

DOI Proposal

Sampling

General Resource Type

Keywords (free)



| | General Resource Type |
|---|---|
| | Resource Type |
| 0 | |
| 0.1 | |
| | Title |
| 1 | Other Titles |
| | Title Type |
| 2 | Collective Title |
| 2.1 | Numbering |
| 3 | Creator |
| 3.1 | |
| 4 | Person |
| | Person ID |
| 4.1 | Vocabulary of Person ID |
| 4.1.1 | URI Name Authority |
| 4.1.1.1 | Record[2] |
| | Affiliation |
| 4.1.1.2 | Affiliation ID |
| 4.1.2 | Vocabulary Affiliation ID |
| 4.1.2.1 | URI Institution Authority |
| 4.1.2.1.1 | Record |
| 4.1.2.1.2 | Institution |
| | Institution ID |
| 4.2 | Vocabulary of Institution |
| 4.2.1 | ID |
| | URI Institution Authority |
| 4.2.1.1 | Record |
| 4.2.1.2 | URL |
| | DOI Proposal |
| 8 | Version |
| 9 | Internal Resource |
| 10 | Identifier |
| 10.1 | |

FOSTER

cessda

# Organizing
# Folders and Files

# Structuring Folders

- Systematically managing folders
  - saves time and effort
  - simplifies the use (collaborative projects)
  - protects your folders and files from accidental clean-up

- Hierarchical structure of folders
  - structure by topic, data type etc.

- Develop standards early in the project
  ⇒ use these standards consistently within a project
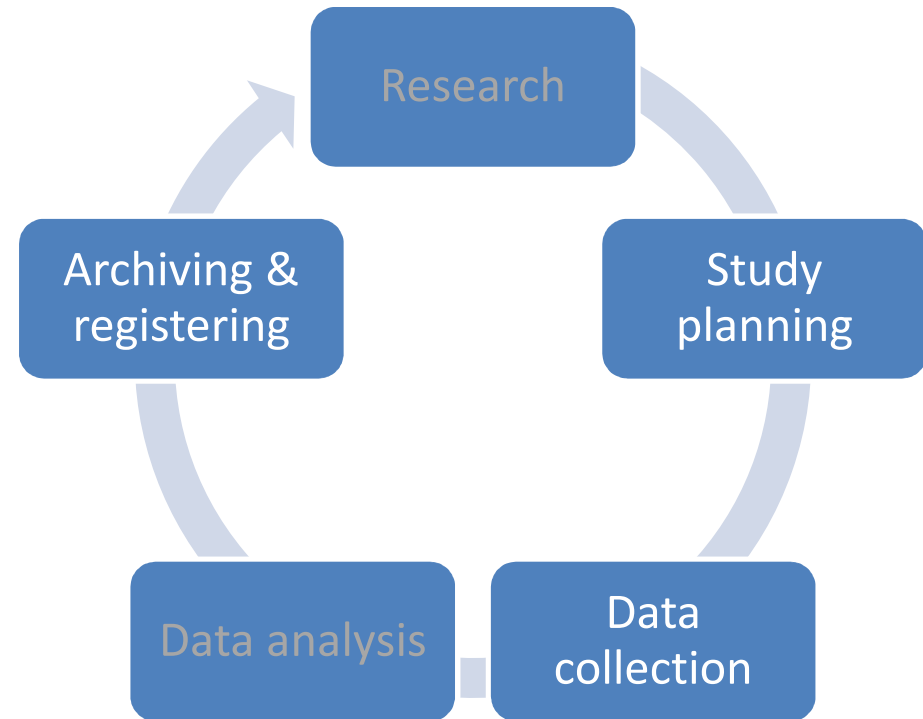
# File Names and Versions

- File names
  - can contain various information, e.g. title of project, editor's name, date of creation, version etc.
  - neither include punctuation characters or blanks nor be too long

- File versioning
  - as a part of the file names, e.g. including the date or numbering the files
  - included in the header of the file
  - in a separate log-file

# (Recommended) File Formats

| Type of data | Recommended formats | Acceptable formats |
|---|---|---|
| Tabular data with extensive metadata | SPSS portable format (.por)<br><br>delimited text and command ('setup') file (SPSS, Stata, SAS, etc.) | SPSS (.sav); Stata (.dta); MS Access (.mdb/.accdb) |
| Tabular data with minimal metadata | comma-separated (.csv); tab-delimited file (.tab) | MS Excel (.xls/.xlsx); MS Access (.mdb/.accdb),<br><br>dBase (.dbf); OpenDocument (.ods) |
| Textual data | Rich Text Format (.rtf); plain text, ASCII (.txt) | HTML(.html); MS Word (.doc/.docx); software-specific formats, e.g. NUD*IST or NVivo |
| Image data | TIFF 6.0 uncompressed (.tif) | JPEG (.jpeg, .jpg); RAW image format (.raw), Photoshop (.psd); PDF/A or PDF (.pdf) |
| Audio data | Free Lossless Audio Codec (.flac) | MPEG-1 (.mp3); Waveform (.wav) |
| Video data | MPEG-4 (.mp4); JPEG 2000 (.mj2) | |
| Documentation and scripts | Rich Text Format (.rtf); PDF/A or PDF (.pdf); HTML (.htm); OpenDocument (.odt) | plain text (.txt); MS Word (.doc/.docx), MS Excel (.xls/.xlsx); XML (.xml) |

Source: UK DATA Service, http://ukdataservice.ac.uk/manage-data/format/recommended-formats

# Data Storage and Security



Research

Study planning

Data collection

Data analysis

Archiving & registering

# Back-Up

- Digital media are fallible

- A back-up is an additional copy that can be used to restore originals

- Backing-up implies having a back-up strategy



ALWAYS MAKE A BACK-UP!!!

A GOOD ADVICE

Image by A. Herrema & H. Bouwteam  (CC-by)

# Towards a Back-up Strategy

- A systematic back-up strategy defines
  a) what          $\Rightarrow$  all, some, just changes …
  b) where         $\Rightarrow$  external, local, remote copies …
  c) how often     $\Rightarrow$  at least in triplicates
  d) for how long  $\Rightarrow$  how long are things needed
  e) responsibility $\Rightarrow$  automate the back-up process

- Verify and recover your back-ups
    $\Rightarrow$ never assume, regularly test a restore

- Treat back-ups the same as the original files

# Data Protection

- Protect your data from unauthorized access, use, change, disclosure, destruction etc.

- Take care of personal data
  - data protection legislation (EU Directive 1995/46/EC)
  - separate personal data from other data

- Use passwords and encryption



Image by P. Hochstenbach (CC-by)

cc-by http://hochstenbach.wordpress.com

# Passwords

- A strong password has
  - eight to fifteen characters or even more
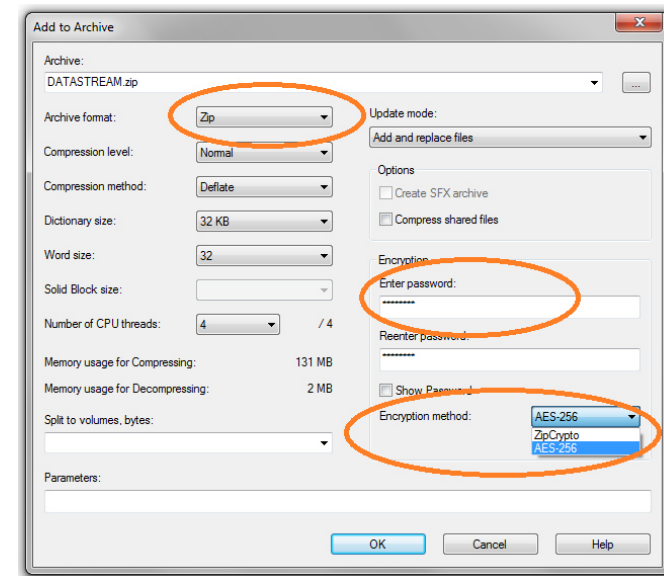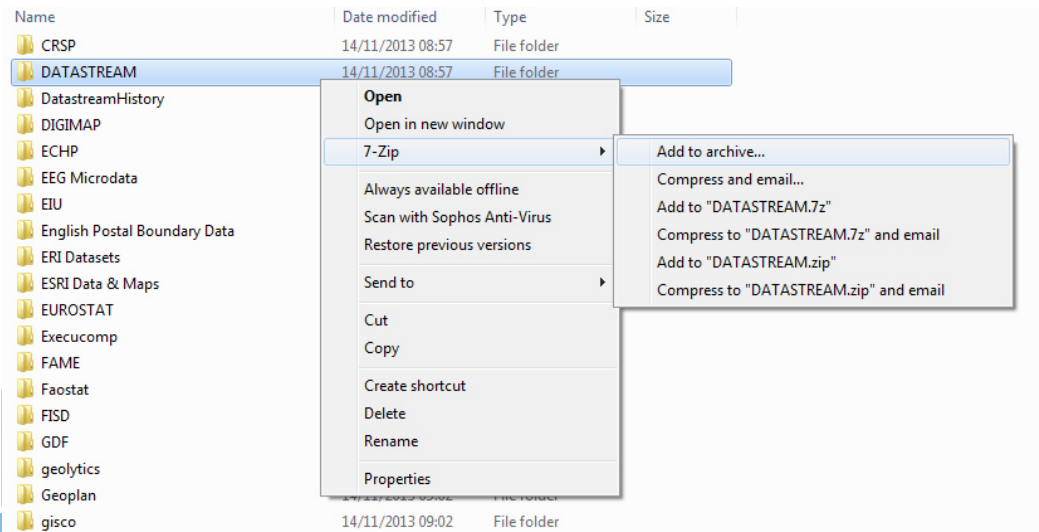  - a random distribution of characters

- Combine…

  
  Image: CC-0

  … upper case letters:   A - Z
  … lower case letters:   a - z
  … numerals:             0 - 9
  … special characters:   ! " # $ % & ' ( ) * + , - . / : *etc.*

# Encryption

- Helps maintain the security of data and documentation
  - uses an algorithm to transform information
  - requires a "key" to decrypt

- For example, encrypt ZIP files securely using 7Zip

# Further Readings

- Aryal, M. (ed.) (2012): Speak Safe. Media Workers' Toolkit for Safer Online and Mobile Practices. https://www.internews.org/sites/default/files/ resources/Internews_SpeakSafeToolkit.pdf

- Borgmann, M., Hahn, T., Herfert, M., Kunz, T., Richter, M., Viebeg, U., and Vowé, S. (2012): On the Security of Cloud Storage Services. Frauenhofer Institut, SIT Technical Report. https://www.sit.fraunhofer.de/fileadmin/dokumente/ studien_und_technical_reports/Cloud-Storage-Security_a4.pdf.

- Directive 95/46/EC of the European Parliament and of the Council, 24 October 1995. Available at: http://eur-lex.europa.eu/LexUriServ/LexUriServ.do? uri=CELEX:31995L0046:EN:NOT

- Gregory, A., Heus, P., & Ryssevik, J., 2009, Metadata. Berlin. http://www.ratswd.de/download/workingpapers2009/57_09.pdf.

- Miller, K., & Vardigan, M., 2005, How Initiative Benefits the Research Community - the Data Documentation Initiative. In First International Conference on e-Social Science, Manchester, UK, June 2005. http://www.ddialliance.org/sites/default/files/miller.pdf.

- National Information Standards Organization, 2004, Understanding Metadata (p. 17). Bethesa, MD: NISO Press. www.niso.org/standards/resources/UnderstandingMetadata.pdf.

- Plant, R. R., 2012, How to add metadata to your data so that you and others can make sense of it. Retrieved from http://www.shef.ac.uk/polopoly_fs/1.158828!/file/Metadatav6.pdf.

- Starr, J., 2011, DataCite Metadata Schema for the Publication and Citation of Research Data (p. 29). doi:10.5438/0005

- Vardigan, M., Heus, P., & Thomas, W., 2008, Data Documentation Initiative: Toward a Standard for the Social Sciences. International Journal of Digital Curation, 3(1), 107–113. doi:10.2218/ijdc.v3i1.45.

# DMP Sections 2 and 3

- 💬 work in 3 groups,

- 🕐 time: about 30 minutes

- ☑ choose one of the following topics

# DMP Sections 2 & 3

a) documentation (*Section 2*), considering …

… what information is needed

… how you (will) capture this information

b) data storage and back-ups (*Section 3.1*), developing a back-up strategy, i.e. …

… what, where and how often / long it is backed-up

… how are back-ups verified

c) managing folders and files (*Section 3.3*), considering …

… how you will organize your folders

… how you will name and version your files

# DMP Section 2: Documentation

- study description
  - study's aim, primary researcher (and funders), population and sampling procedures, method of data collection, data cleaning and anonymization etc.
    - ⇒ technical and methodological report

- variable description
  - questionnaire: original question wording and provided answer categories, explanations, interviewer instructions;
  - variables: labels and meanings of variables and codes, variable notes, scales etc.
    - ⇒ codebook, questionnaire and labels in the dataset

# DMP Section 3.1: Back-up

- developing a back-up strategy
    - ⇒ defining clear and consistent guidelines
        - – what: all, something, only changed files
        - – where: at least in triplicates and different locations
        - – how long are different files (and versions) needed
          never destruct or overwrite original data
        - – who: name researcher(s) and assign responsibilities
    - ⇒ verify back-ups frequently (e.g. once a week),
      e.g. restoring the files (name researcher(s) and
      responsibilities)

# DMP Sections 3.3: Organizing …

- developing guidelines to organize …

  … folders

  - define a consistent structure of folders

    ⇒ e.g. by topic

  … and files, i.e. define a consistent strategy

  - to name files

    ⇒ e.g.  [type_name_version]

  - to version files

    ⇒ e.g. by the date and editor's acronyms

  - *data_RDMData_20150822sn*

RDM-Project
- 1_Organisational
  - Administrative
  - Communications
- 2_Literature
  - Abstracts
  - Articles
- 3_Design
  - Instrument -1
  - Instrument- 2
- 4_Data
  - 1_PrimaryData
    - 1_DataDocumentation
    - 2_Syntax
    - 3_Output
  - 2_SecondaryData
- 5_Presentations
  - Conference-1
  - Conference-2
- 6_Papers
  - 1_Publications
    - Article -1
      - 1_Paper
      - 2_Analyses
      - 3_Reviews
    - Article- 2
  - 2_Proposals
    - 1_Outline
    - 2_Drafts
    - 3_Reviews
- 7_Others