

washr an R-package to facilitate FAIR data publishing

FAIR Data Publishing

Lars Schöbitz

Nov 28, 2024

Data

 hands up 

Who brought data?

 hands up 

Who has got “A one line description of the data”?



Who has got “A one paragraph, three sentence description of the data”?



Who has got a “Data dictionary, stored as dictionary.xlsx. Two columns, describing each variable in the dataset you come with.”

Metadata

Metadata: data about data

WHAT!?

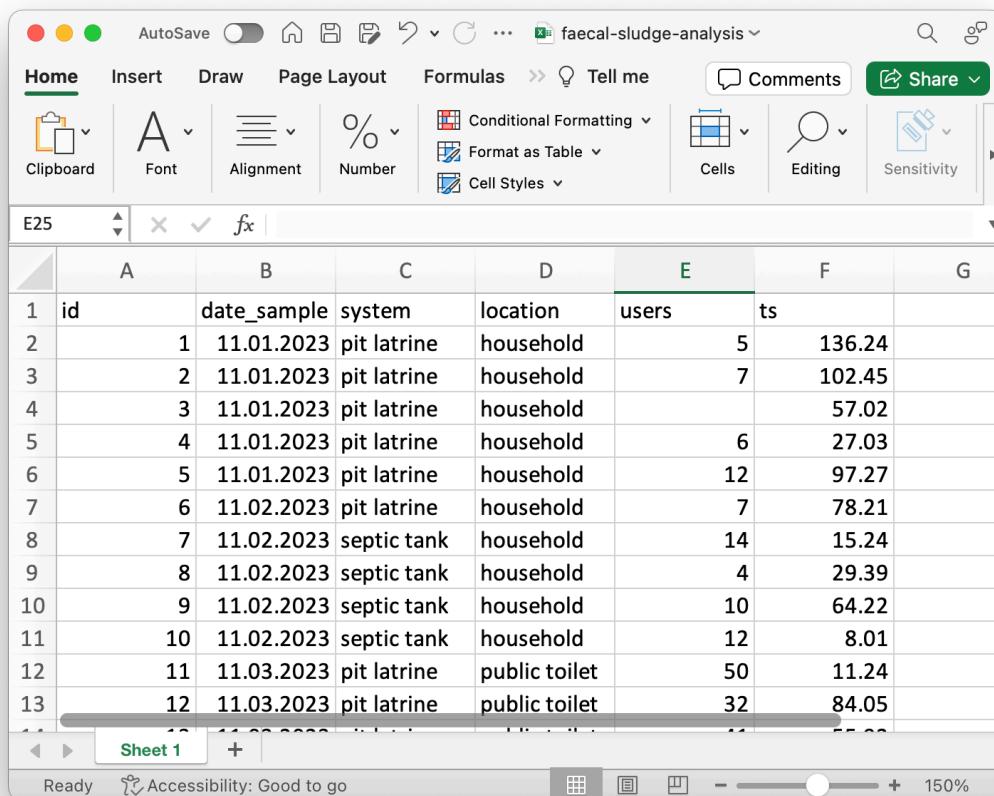
Faecal sludge samples

Imagine:

- you are new to WASH research and you have never heard of faecal sludge management.
- you are interested in learning more about the topic and you want to find some data to play with.
- you find a publication with a dataset on faecal sludge characteristics.

Faecal sludge samples

You download the XLSX file that contains the data and you open it in Excel. You see the following:



The screenshot shows a Microsoft Excel spreadsheet titled "faecal-sludge-analysis". The data is presented in a table with the following columns: id, date_sample, system, location, users, and ts. The rows contain 13 data points, each with a unique ID from 1 to 13, a date between November 2023 and March 2023, a system type (pit latrine or septic tank), a location (household or public toilet), a number of users, and a corresponding value for ts.

	A	B	C	D	E	F	G
1	id	date_sample	system	location	users	ts	
2	1	11.01.2023	pit latrine	household	5	136.24	
3	2	11.01.2023	pit latrine	household	7	102.45	
4	3	11.01.2023	pit latrine	household		57.02	
5	4	11.01.2023	pit latrine	household	6	27.03	
6	5	11.01.2023	pit latrine	household	12	97.27	
7	6	11.02.2023	pit latrine	household	7	78.21	
8	7	11.02.2023	septic tank	household	14	15.24	
9	8	11.02.2023	septic tank	household	4	29.39	
10	9	11.02.2023	septic tank	household	10	64.22	
11	10	11.02.2023	septic tank	household	12	8.01	
12	11	11.03.2023	pit latrine	public toilet	50	11.24	
13	12	11.03.2023	pit latrine	public toilet	32	84.05	

Faecal sludge samples

Open questions:

- What unit does **users** refer to?
- What does **ts** stand for?
- The **date** of what?
- Where was this data collected?
- Which method was used to collect the samples?

The screenshot shows a Microsoft Excel spreadsheet titled "faecal-sludge-analysis". The table has columns labeled A through G. Column A contains the index number (1 to 13). Columns B, C, and D contain categorical data: date_sample, system, and location respectively. Columns E and F contain numerical data: users and ts. The data shows 13 rows of information, with the last row partially visible.

	A	B	C	D	E	F	G
1	id	date_sample	system	location	users	ts	
2	1	11.01.2023	pit latrine	household	5	136.24	
3	2	11.01.2023	pit latrine	household	7	102.45	
4	3	11.01.2023	pit latrine	household		57.02	
5	4	11.01.2023	pit latrine	household	6	27.03	
6	5	11.01.2023	pit latrine	household	12	97.27	
7	6	11.02.2023	pit latrine	household	7	78.21	
8	7	11.02.2023	septic tank	household	14	15.24	
9	8	11.02.2023	septic tank	household	4	29.39	
10	9	11.02.2023	septic tank	household	10	64.22	
11	10	11.02.2023	septic tank	household	12	8.01	
12	11	11.03.2023	pit latrine	public toilet	50	11.24	
13	12	11.03.2023	pit latrine	public toilet	32	84.05	

Questions that only the original author may have the answers to.

You as an author

have the chance to document your data properly once to make it easier for everyone else to know what it contains.

Documentation

Goes into a separate README file

- General information (authors, title, date, geographic location, etc.)
- Sharing / access information (license, links to publications, citation)
- Methodological information (sampling, analysis, etc.)

Data dictionary

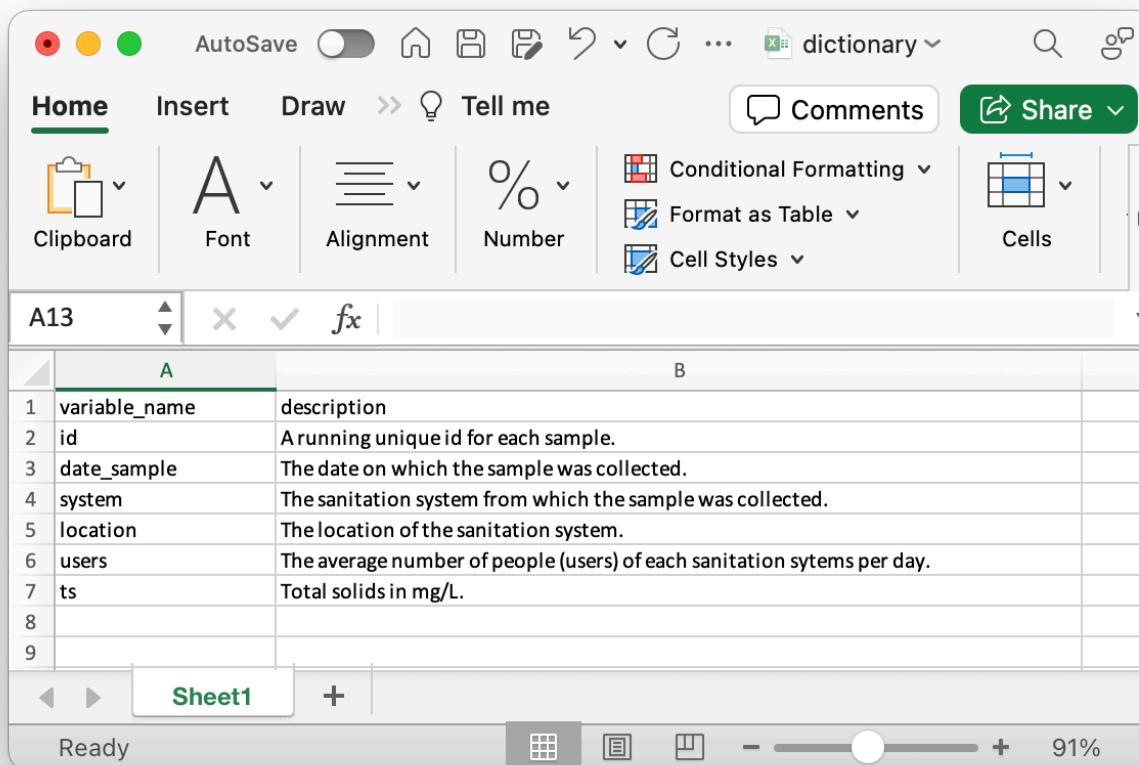
Goes into a separate file (`dictionary.csv`).

Minimum required information

- Variable name
- Variable description

Data dictionary for faecal sludge samples

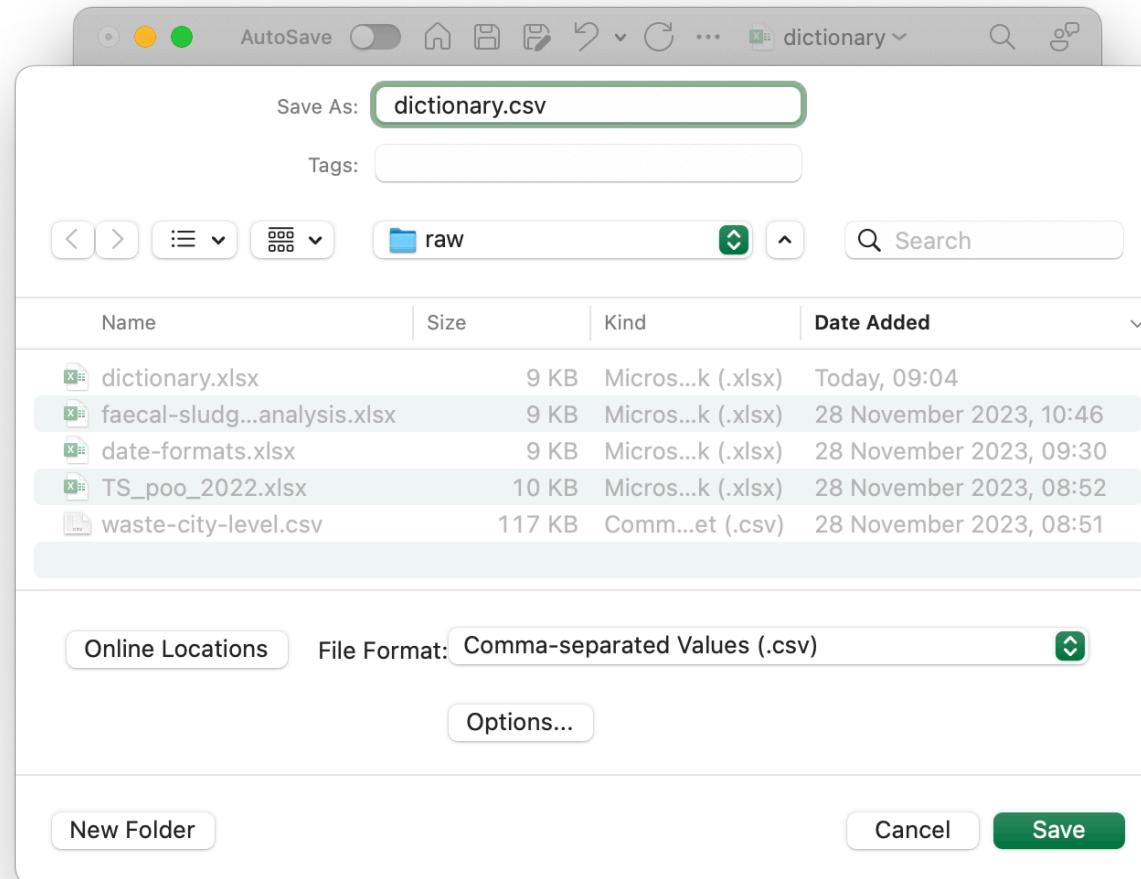
- Edit in spreadsheet software (e.g. MS Excel)



The screenshot shows a Microsoft Excel spreadsheet titled "dictionary". The ribbon is visible at the top with tabs for Home, Insert, Draw, Tell me, Comments, and Share. The Home tab is selected. The formula bar shows "A13". The main content is a table with two columns, A and B. Column A contains variable names and column B contains their descriptions.

	A	B
1	variable_name	description
2	id	A running unique id for each sample.
3	date_sample	The date on which the sample was collected.
4	system	The sanitation system from which the sample was collected.
5	location	The location of the sanitation system.
6	users	The average number of people (users) of each sanitation systems per day.
7	ts	Total solids in mg/L.
8		
9		

Data dictionary for faecal sludge samples



- Save as CSV file

Directory tree of an R data package

Capstone project of Rainbow Train: <https://github.com/openwashdata/washopenresearch> ↗

```
•
├── CITATION.cff
├── DESCRIPTION
├── LICENSE.md
├── NAMESPACE
└── R
    ├── README.Rmd
    ├── README.md
    └── _pkgdown.yml
├── data
├── data-raw
├── docs
├── inst
├── man
└── vignettes
washopenresearch.Rproj
```

Directory tree of an R data package

- `R` folder: R scripts with documentation of data resources
- `data-raw` folder: raw data, contains `data_processing.R` script to prepare processed, analysis-ready data
- `data` folder: documented analysis-ready data in R-internal `.rda` format
- `README.Rmd`: R data package documentation, displayed on GitHub and website landing page

Inside the data-raw folder

```
data-raw/
└── data_processing.R
    ├── dictionary.csv
    └── unc-article-url-manual-
        collection.csv
            washdev.csv
```

- `unc-article-url-manual-collection.csv`: one raw data resource
- `washdev.csv`: another raw data resource
- `data_processing.R`: where the work happens
- `dictionary.csv`: data dictionary describing the variables in the processed, analysis-ready data

My turn: A tour of washopenresearch

Sit back and enjoy!

Your turn: Create a GitHub repository for your data

1. Navigate to the GitHub organisation for this course:

github.com/fairdatapub-washcentre/ ↗

2. Click on the green “New” button to create a new repository.

3. Choose a name for your repository (all small letters, no spaces or dashes)

4. Make the repository public.

5. Keep everything else unchanged.

6. Click on green “Create repository” button.

7. Keep the tab open for the next steps.

Our turn: Create a project on Posit Cloud

1. Navigate to the course workspace on [posit.cloud/spaces/588816/content/↗](https://posit.cloud/spaces/588816/content/).
2. If you haven't yet, bookmark the page.
3. Follow my instructions on the screen.

Take a break

Enjoy your lunch! Let your emails rest in peace.



Our turn: Work through Data Publishing Guide

For the rest of the day
we will work through

Your turn: Prepare a dataset on your own

1. Navigate to the course website: fairdatapub-washcentre.github.io/website/
2. In the left-hand menu, click on Module 3, then select am-03:
Data Organization
3. Follow the instructions
4. Place a yellow sticky note on your laptop when you have completed the assignment

Module 3 documentation

fairdatapub-washcentre.github.io/website/modules/md-03.html

Wrap-up

Participate in our hackathon tomorrow

- All tools welcome
- Win up to 750 CHF (~ 850 USD)
- <https://openwashdata.org/pages/events/2024-11-18-hackathon/>
- Registration is still open: https://forms.office.com/pages/responsepage.aspx?id=7KY0lmaioWrFHTEIR_FgjqEYzwAXr5PiYenM5_ZB_lURTU2UklP

Your turn: Sign up for the openwashdata newsletter!

 fairdatapub-washcentre.github.io/website/

05 : 00

What's next?

Holiday break! I will be off from December 9th to January 15th.



What's next?

Happy New Year! I will be at UKZN on Thursday, 16th January. I will leave South Africa on Thursday, 13th February.



Thanks! 🌻

Slides created via revealjs and Quarto:

<https://quarto.org/docs/presentations/revealjs/>

Access slides as [PDF on GitHub](#) ↗

All material is licensed under [Creative Commons Attribution Share Alike 4.0 International](#) ↗.