# Agroecology Partnership: Data Management Guidelines

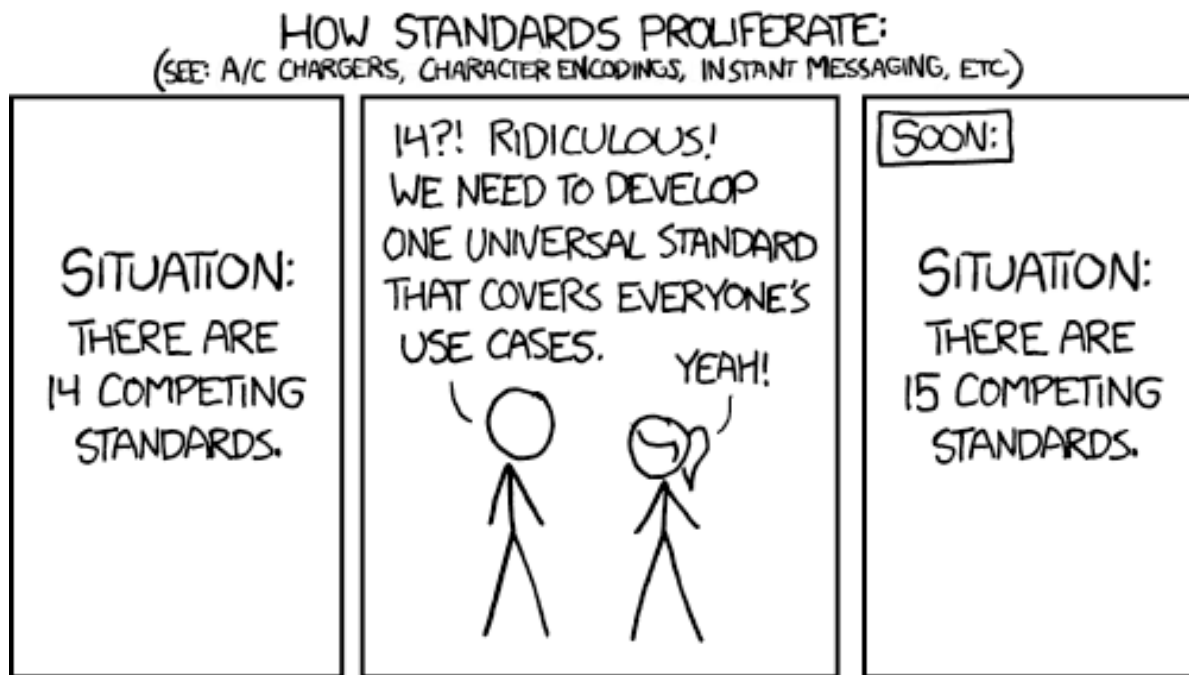LifeWatch ERIC

August 5, 2025

# Table of contents

# Introduction

Welcome to the Data Management Guidelines documentation of the Agroecology Partnership.. This online documentation expands on the official Agroecology Partnership Data Management Plan (DMP) and complements it with more practical and technically detailed information.

While the DMP outlines the full strategy for data handling across the partnership and it is updated every 2 years, the present documentation is meant to be updated as we learn from doing. Its objective is to support partners with hands-on guidance, concrete examples, and step-by-step explanations of the data pipeline — from collection to publication — including how to make data FAIR (Findable, Accessible, Interoperable, and Reusable). This space aims to be more to the point than the DMP, helping teams implement the plan more effectively in day-to-day work. Furthermore, being a public document, we aim to offer the wider agronomy community with state-of-the-art data management practices in the Agroecology domain.

For any questions, please write us a line at agroecology-data@lifewatch.eu.

# Data Standarization

Data on Agroecology are multidisciplinary: There are biological data, socio-economic, geographical and many more. The main goal is to standardize data to the preferred community standard. However, sometimes there is not a clear winning standard, or the winning standard may not meet the FAIR principles. This chapter describes the preferred data standards for each data type collected in the Agroecology Parnership. In some cases, when there is not an obvious standard, we suggest evidence-based way of structuring data that fits the FAIR principles.



Source: xkcd

## Survey data

Recommendations for organizing survey data in the Agroecology partnership following a wide, spreadsheet format approach. This structure is inspired by the data structure of the American National Election Studies (ANES 2020), the recommendations of (Zimmer, Powell, and

Velásquez 2024), while applying a tidy data approach (Wickham 2014) with usability on mind and the enabling of a later transformation into Open Linked Data as structured in The Survey Ontology (Scrocca et al. 2021).

## Survey project structure

We propose a project structure that uses `csv` files as a way of formally describing the survey, both questions, answer and other information. This structure may be accompanied of a text description of the survey for further context. A minimal example of the file directory would look like:

```
.
    data/
      codebook.csv
      responses.csv
    docs/
       survey_descriptor.pdf
```

The data files are inside the data directory, and the different documents including the description of the survey are inside the docs directory.

This structure can be further extended, depending on the needs of the survey. For instance, a README file can be added to the root of the project with a short explanation of the survey. When uploading to Zenodo, this can be written instead in the Description of the metadata. If the survey is uploaded to GitHub, we recommend to use the markdown format for the README file. Also on GitHub, it is recommended to include a LICENSE file, choosing the license of the data.

In addition, any notebook or spreadsheet used to analyze the data can be added to the analysis folder, and reusable code scripts can be added to the scripts directory.

```
.
  README.md
  LICENSE.md
    data/
      codebook.csv
      responses.csv
    docs/
      survey_descriptor.pdf
      survey_analysis.pdf
    scripts/
      clean_data.R
```

```
analysis/
    survey_analysis.xlsx or survey_analysis.ipynb
```

**Survey project description file**

A survey is always quite specific, and cannot be fully understood without a good explanation
of what is the context and the rationale behind the different questions. It is also a sign of
respect to the respondents to explain what is the goal of your survey. We propose to include
this description in a `survey_descriptor.pdf` file in the `docs/` directory.

**Survey project data files**

Surveys are typically are typically stored as `pdf`, `docx` or `xlsx`. Having interoperability on
mind, we propose to use **text delimited files** such as `csv`. These type of files are largely used
in data science. they have several advantages, including easy machine-readability and being
an open format with no owner, which ensures data will remain readable and understandable
by many different software for a long time. We propose the following specifications:

- Separate columns using semicolons ;
- Use double-quotes " to quote strings
- Use `UTF-8` as encoding.

We avoid using colons , as separators because these can be used in free text, open questions.
They would affect the structure of the data. We recommend to prohibit the use of semicolons
; and double-quotes " in any case that is not separating columns and delimiting strings
respectively. We encourage you to use software solutions that allow blocking these characters
in free-text answers provided by users. If blocking the use of these characters it is not possible,
we urge you to ask respondents to not use them while filling up your survey.

**Codebook**

Codebooks are files that **explain the questions** formulated in the survey. An unique identifier
(a code) is assigned to each question, linking the information about the questions with the
answers provided by the participants. But the codebook can be used as well during the design
of the survey.

The example below shows an hypothetical survey codebook about the user satisfaction using
an online platform.

`./data/codebook.csv`

| Code | Label | Type | QuestionText | Values | Cardinality | QuestionType |
|------|-------|------|--------------|--------|-------------|--------------|
| Q1_age | Age | integer | What is your age? | | 1..1 | SingleChoice |
| Q2_gender | Gender | string | What is your gender? | Female|Male|Other | 1..1 | SingleChoice |
| Q3_satisfaction | Satisfaction | integer | How satisfied are you with the platform? | 1=Very unsatisfied|2=Unsatisfied|3=Neutral|4=Satisfied|5=Very satisfied | 1..1 | SingleChoice |
| Q4_improvement | Improvement | string | What would you improve in the platform? | | 0..n | FreeText |

Each row in the codebook describes a survey question. Below is an explanation of each column:

- **Code**: A unique identifier for the question. It must be unique within the survey.

- **Label**: A short, human-readable label or name for the variable that can be used in spreadsheets or statistical software.

- **Type**: The data type of the answer. Common types include "integer", "string", "boolean", or "date".

- **QuestionText**: The full text of the question as it was asked in the survey.

- **Values**: A list of possible values for closed questions. Options are separated by vertical bars (|), and value labels can be assigned using the equals sign (=). For open or free-text questions, this field is left empty.

- **Cardinality**: Indicates how many answers are allowed. The cardinality pattern is `min..max`, where the first number is the minimum required answers and the second is the maximum allowed. `n` denotes "no fixed upper limit," so `1..1` means exactly one answer, while `0..n` means the question may be skipped or answered multiple times.

- **QuestionType**: Describes the nature of the question. Typical values include "Single-Choice", "MultipleChoice" or "FreeText"

## Responses

**Answers to the survey** are recorded in the `./data/responses.csv` file. **Every row is the answers of a participant**, and **every column is named after the `code`** in `codebook.csv`. This allows to link easily the information about the questions without getting the responses file full of details that difficult the analysis.

`./data/responses.csv`

| respondent_id | Q1_age | Q2_gender | Q3_satisfaction | Q4_improvement |
|---|---|---|---|---|
| 001 | 34 | female | 4 | I think the platform is user-friendly |
| 002 | 29 | male | 5 | Needs better support for collaboration. |

**We recommend to include an unique identifier for every respondent**, here named as `respondent_id`. This allows to annomymize the survey without loosing the link to private information about the respondents that might have been collected (e.g. name, email, address)

## Participants and private information

Sometimes your survey might collect personal or sensitive information — like names, emails, or locations. If that's the case, you **should never share this data publicly**. In the Agroecology Partnership, we follow the **GDPR** (the General Data Protection Regulation), which means that **private data must be handled carefully and securely**.

**We recommend creating a separate file** (not to be published!) called `participants.csv` inside your `data/` folder. This file is **for internal use only** — for example, if you need to follow up with participants or validate something later. To keep things tidy and safe, give each person a unique code, which we'll call `respondent_id` and store personal info (like name or email) in separate columns. Use clear, human-friendly column names, and write them using `snake_case` (i.e. separate the words of the column names with underscores `_`).

| respondent_id | name | email | country | preferred_language | wants_emails |
|---|---|---|---|---|---|
| 001 | Jane Doe | jane.doe@example.com | France | fr | FALSE |
| 002 | John Doe | john.doe@example.com | USA | en | TRUE |

The table above is an example (with mocked information) of how the `participats.csv` would look like. This file helps you **anonymize your actual survey data**, because in the main `responses.csv`, you'll only keep the `respondent_id`, and this is nothing that could personally identify someone.

# References

ANES. 2020. "Time Series Study Full Release: User Guide and Codebook." https://electionstudies.org/wp-content/uploads/2022/02/anes_timeseries_2020_userguidecodebook_20220210.pdf.

Scrocca, Marco, Daniele Scandolari, Giulia Re Calegari, Irene Baroni, and Irene Celino. 2021. "The Survey Ontology: Packaging Survey Research as Research Objects." In *Proceedings of the 2nd Workshop on Data and Research Objects Management for Linked Open Science – Co-Located with ISWC 2021.* https://doi.org/10.4126/FRL01-006429412.

Wickham, Hadley. 2014. "Tidy Data." *Journal of Statistical Software* 59 (10): 1–23. https://doi.org/10.18637/jss.v059.i10.

Zimmer, Samantha A., Ryan J. Powell, and Iván C. Velásquez. 2024. *Exploring Complex Survey Data Analysis Using r: A Tidy Introduction with {Srvyr} and {Survey}.* Chapman & Hall/CRC Press. https://tidy-survey-r.github.io/tidy-survey-book/.