# A brief introduction on data curation process

Pietro Bartolini

IFAD EC Drylands Restoration Team Meeting
6-8 May, Nairobi, Kenya

ICRISAT
INTERNATIONAL CROPS RESEARCH
INSTITUTE FOR THE SEMI-ARID TROPICS

World Agroforestry 40

ILRI
INTERNATIONAL
LIVESTOCK RESEARCH
INSTITUTE

IFAD
Investing in rural people

The Dataverse Project

ICARDA
Science for resilient livelihoods in dry areas

**icarda.org**
International Center for Agricultural Research in the Dry Areas

cgiar.org
A CGIAR Research Center

CGIAR

# Introduction

A dataset, in various form, is a by-product of many research activities. Unfortunately, data curation is often neglected in comparison to the writing of an article or a report.

During the research, it is common to organize data in spreadsheets in a way which makes them easily understandable for the author at that specific time.

However the presentation of the data is very important to improve the spreading of the research, even if the dataset is not the primary product.

If the dataset is going to be released to the public, we must follow machine-readable standards, assuring that any future user can read it, understand its content, and use the data in other research projects.

# Preliminary Steps

"When you are working with spreadsheets, during data clean up or analyses, it's very easy to end up with a spreadsheet that looks very different from the one you started with" (Bahlai, 2017).

In order to avoid errors and data losses, do not modify the original dataset. Create a new copy to curate.

Enhance the title of the new dataset, providing some context.

| Old Dataset Title: | Rangeland Species Composition |
|---|---|
| Enhanced Dataset Title: | Annual and Perennial Rangeland Plant Cover and Species Composition, Tatatouine, Tunisia, November 2018 |

If the data is from a Primary Article Citation, use the naming convention "Data from: title of the article" (USDA, 2016).

# Data Curation – Elaboration Management

The dataset should contain only raw data. Each elaboration is subject to error. Providing only raw data, we enable future users to use them for their research, avoiding the risk to replicate elaboration errors.

- No graphs

- No formulas

- No percentages

- No elaborations

# Data Curation – Tables Arrangement

It is not possible to have multiple data tables in one spreadsheet and use blank rows or columns to separate the data. It is not machine-readable.

Each spreadsheet must contain a single table.

# Data Curation – Tables Arrangement

If different spreadsheets contain similar data, that need to be analysed together, they can be merged, allowing the computer to see the connection.

However, be sure to add a column defining possible differences, like date or location.

# Data Curation – Data Management

The dataset structure should be clean and simple:

- Use short title for the head columns, without spaces or symbols. Just write in camel case or use underscore

- Only 1 information for each cell

- No vague or misleading information, when possible use numeric code

- Use ISO standards for  date and time (YYYYMMDDHHMM)

- No merged cells

- No comments, use a column for the notes

- No empty cells. Use NA for missing or null data

- Although text format and colour could improve the readability they must be avoided, because they can easily be lost during transfer, compromizing the overall structure

# File Format

Whatever software we are using we must be sure the files will be readable in the future, using different version of a licensed or unlicensed software.

The CSV or comma separated value files are the preferred data format for most of data repositories and are the recommended one for publishing machine-readable tabular data.

A CSV file contains a single spreadsheet: a dataset uploaded in MEL will be a collection of several CSV files.

Files and Links:

✗ Data_Introduction.csv 📥⦿ Mark as main file ❶

✗ Data_Element_Description.csv 📥○ Mark as main file ❶

✗ Unique_Identifier.csv 📥○ Mark as main file ❶

✗ Dams.csv 📥○ Mark as main file ❶

✗ Tanks.csv 📥○ Mark as main file ❶

✗ Check_Dams.csv 📥○ Mark as main file ❶

✗ Contour_Structures.csv 📥○ Mark as main file ❶

✗ Reforestation.csv 📥○ Mark as main file ❶

✗ Desert_Restoration.csv 📥○ Mark as main file ❶

✗ Spontaneous_Intervention.csv 📥○ Mark as main file ❶

# Data Dictionary – Dataset Introduction

The Dataset Introduction provide an overall explanation about the dataset scope and creation. It must include:

- Description: A rich free text description that provides as much explanation as possible about the dataset: how and why it was generated, and how it should (or should not) be used. Make sure that in this description are present the experiment settings (location, climatic conditions, etc.), data collection and processes methods, equipment used, period, possible resources and any limiting factors (USDA, 2016)

- Summary: A shorter description of the dataset, usually no more than a sentence or two (USDA, 2016)

- Start_Date: The date in which the data collection starts

- End_Date: The date in which the data collection ends

- Author: Dataset first author

- CoAuthor: Dataset co-authors

# Data Dictionary – Dataset Introduction

| Description | Summary | Start_Date | End_Date | Author | CoAuthor |
|---|---|---|---|---|---|
| A rich free text description that provides as much explanation as possible about the dataset. | A shorter description of the dataset, usually no more than a sentence or two. | YYYYMMDD | YYYYMMDD | Dataset first author | Dataset co-authors |

The tab structure is customizable according to the needs of the author. It can include specific section about geographical location, metodologies and additional notes.

# Data Dictionary – Element Description

This is the most important component of the Data Dictionary. It provides explanation about the meaning of each variable and correspondences for any code used.

| Spreadsheet_Tab | Element_DisplayName | Description | Units | Data_Type | Character_Length | Acceptable_Values | Required | Accepts_NullValue |
|---|---|---|---|---|---|---|---|---|
| Spreadsheet_Name | Spreadsheet_Name | Description of the spreadsheet content | NA | NA | NA | NA | NA | NA |
| Spreadsheet_Name | Variable_N1 | Description of the variable meaning | Kg | Numeric | 255 | [x, z] | Y/N | Y/N |
| Spreadsheet_Name | Variable_N2 | Description of the variable meaning | NA | Numeric | 2 | x\|y\|z | Y/N | Y/N |
| Spreadsheet_Name | Variable_N3 | Description of the variable meaning | NA | Text | 255 | NA | Y/N | Y/N |
| Spreadsheet_Name | Variable_N3 | Description of the variable meaning | YYYYMMDD | Date | 8 | [yyyymmdd, YYYYMMDD] | Y/N | Y/N |

# Data Dictionary – Element Description

The suggested template for structuring manually the "Dataset Elements Description" includes the following fields (USDA, 2016):

- Spreadsheet_Tab: The tab where is the element

- Element_DisplayName: The dataset element name

- Description: A brief and complete element definition that could stand alone from other elements definition

| B | C |
|---|---|
| Element_DisplayName | Description |
| number | Invoice number |
| date | Invoice date |
| status | Invoice status |
| amount | Invoice amount |
| customer_no | Customer number |

| B | C |
|---|---|
| Element_DisplayName | Description |
| number | Invoice autogenerated number, starting from 1 each year. Number is generated when invoice gets approved. |
| date | Invoice issued date. Null for working copy invoices. Automatically set to today's date on invoice approval. |
| status | Invoice status. 'W' - working copy, 'A' - approved invoice, 'C' - cancelled. |
| amount | Invoice net amount in USD |
| customer_no | Number of customer invoice was issued to. Ref: customers. |

# Data Dictionary – Element Description

- Unit: The unit of measurement adopted for the elements
- Data_Type: The type of data values contained in the field (e.g. varchar, integer, date, etc.)
- Character_Length: The length of data values contained in the field (maximum length for Excel is 255)
- Acceptable_Values: The list of acceptable values in this field. In some case it can be also a range of values
- Required: Express the requirement of values in the field for dataset status and validity
- Accepts_NullValue: Express the possibility of null values in the corresponding dataset field



If the dataset includes multiple tabs, the various element of each tab should be listed in the same Data Dictionary file.

# Data Dictionary – Unique Identifier

We often assume the use of certain terms to be clear, but it is not always the case, especially outside our research team.

To make sure to solve any possible ambiguity, in the unique identifier tab are reported the corresponding link for the dataset terms and concepts to the on-line thesaurus. This is very useful to avoid any misunderstanding on the elements (plant species, animals, etc.) analyzed and reported in the set of data (Bonechi 2018).

| Spreadsheet_Tab | Element_DisplayName | Unique_Identifier | Source |
|---|---|---|---|
| Spreadsheet_Name1 | Earth dams | http://aims.fao.org/aos/agrovoc/c_32435 | AGROVOC |
| Spreadsheet_Name2 | Sheep fattening | http://lod.nal.usda.gov/nalt/92111 | USDA |
| Spreadsheet_Name3 | Barley | http://aims.fao.org/aos/agrovoc/c_823 | AGROVOC |
| Spreadsheet_Name3 | Malting Barley | http://aims.fao.org/aos/agrovoc/c_25485 | AGROVOC |

# Final Recommendation

http://repo.mel.cgiar.org/handle/20.500.11766/9400

The General Dataset Curation Guide is available as a final draft at the link above. Each one of you can read it and start applying it to his own work until it will become common practice.

The Guide defines a standard, however it is still a work in progress, and certain types of file may need more curation work than others or even *ad hoc* solution. You can signal potential issues to the data curation staff, helping to expand the future version of the guide.

# Thanks for your attention!