







General Dataset Curation Guide

Francesco Bonechi







Introduction

• "Many valuable datasets are poorly curated, which contributes to errors, redundant effort, and obstacles to replication and use" (Ruggles, 2018)

• It is common to organize data in spreadsheets in a way which makes them easily understandable for the dataset author at that time, without following the machine-readable standards or considering any next research use.

Data Curation Role

"Data curation activities enable data discovery and retrieval, maintain data quality, add value, and provide for reuse over time" (Munoz, 2017)

Nowadays specific jobs related to data curation responsibilities are increasing (with title like "data curator" or "data curation specialist")



Preliminary Steps

- 1. In order to be able to reproduce your analyses and be sure you don't lose any data within curation processes, don't modify the original dataset but first create a new copy of your original data where to work on.
- 2. Enhance the dataset title of the new file "entering a descriptive title including dates, locations, and specific metrics that make the dataset unique" (USDA, 2016).

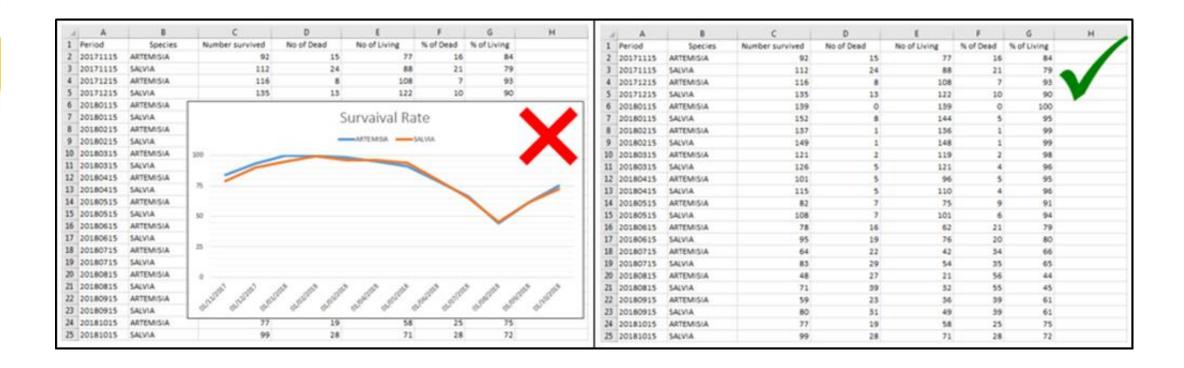
Old Dataset Title: Rangeland Species Composition

Enhanced Dataset Title: Annual and Perennial Rangeland PlantCover and SpeciesComposition Pavlodar

Kazakhstan_November2018

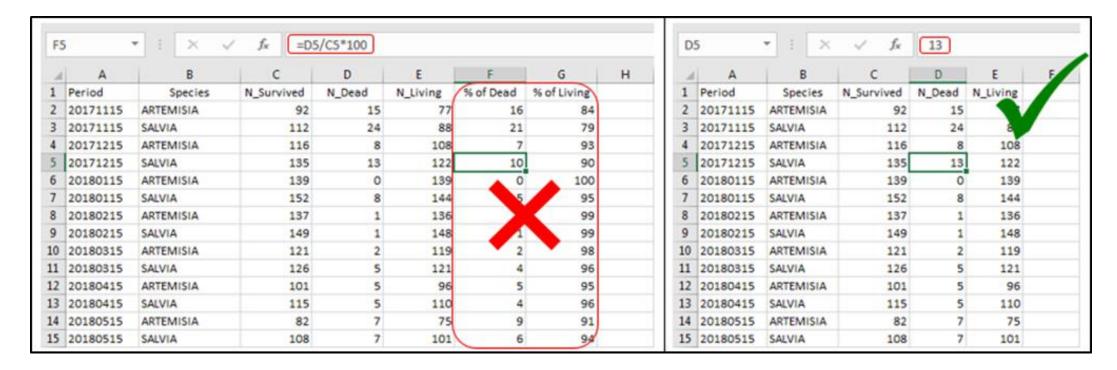
Data Calculations and Summaries

1. Graphics and figures must be removed from the dataset spreadsheet tab.



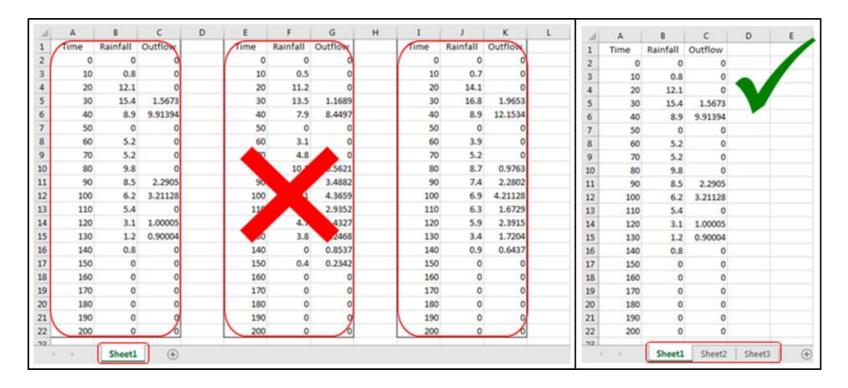
Data Calculations and Summaries

2. Formulas and any other type of elaborations must be removed from the dataset spreadsheet tab. This may force to delete entire spreadsheet columns.



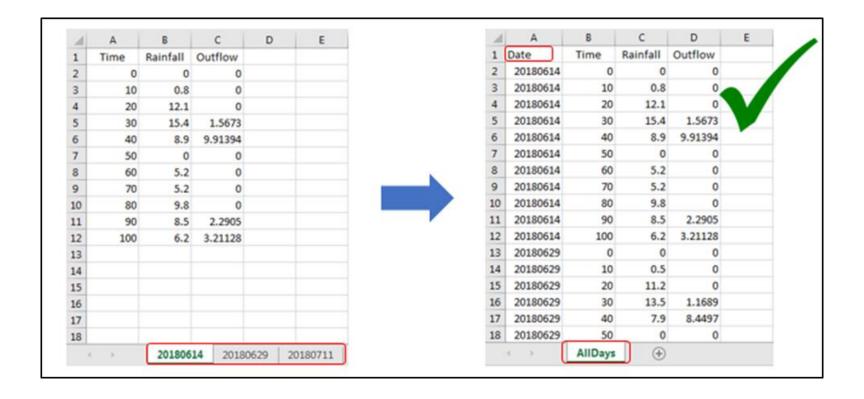
Data Table

1. In order to avoid false association during data system readability, blank rows and columns must not be used to separate the dataset in different tables or sections.



Data Table

Note: Improving the dataset column arrangements is possible to reduce the spreadsheet tab number.



"The cardinal rules of using spreadsheet programs for data are" (Bahlai, 2017):

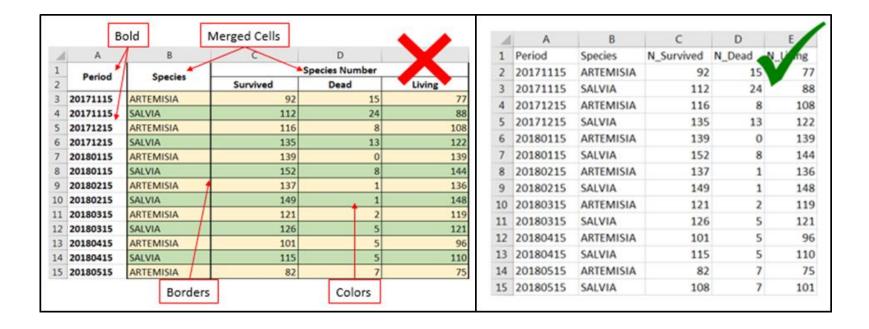
- Put all the variables in columns. Each column corresponds to a variable.
- Put each observation in its own row. Each row corresponds to an observation.
- Don't combine multiple pieces of information in one cell. Each cell corresponds to one value (data).

2 Wheat 678 663.3 548.4 596.1 540 3 Durum 128.8 126.3 117.1 124.8 110.7 4 Barley 647.5 625.2 463.7 574.7 513.2 5 Triticale 10.5 14.1 16.9 15 24.5 6 2015 540 110.7 513.2 6 2015 540 110.7 513.2 24.													
1 Year 2011 2012 2013 2014 2015 2 Wheat 678 663.3 548.4 596.1 540 3 Durum 128.8 126.3 117.1 124.8 110.7 4 Barley 647.5 625.2 463.7 574.7 513.2 5 Triticale 10.5 14.1 16.9 15 24.5 6 2015 540 110.7 513.2 6 2015 540 110.7 513.2 6 2015 540 110.7 513.2 2 2011 678 128.8 647.5 10. 3 2012 663.3 126.3 625.2 14. 4 2013 548.4 117.1 463.7 16. 5 2014 596.1 124.8 574.7 1 6 2015 540 110.7 513.2 24.		Λ	D	C	0	С.	Е	1	Α	В	С	D	E
2 Wheat 678 663.3 548.4 596.1 540 3 Durum 128.8 126.3 117.1 124.8 110.7 4 Barley 647.5 625.2 463.7 574.7 513.2 5 Triticale 10.5 14.1 16.9 15 24.5 6 2015 540 110.7 513.2 6 2015 540 110.7 513.2 6 2015 540 110.7 513.2 2 2011 678 128.8 647.5 10. 3 2012 663.3 126.3 625.2 14. 4 2013 548.4 117.1 463.7 16. 5 2014 596.1 124.8 574.7 1 6 2015 540 110.7 513.2 24.	-4						F	1	Year	Wheat	Durum	Barley	Triticale
2 Wheat 678 663.3 548.4 596.1 540 3 Durum 128.8 126.3 117.1 124.8 110.7 4 Barley 647.5 625.2 463.7 574.7 513.2 5 Triticale 10.5 14.1 16.9 15 24.5 6 2015 540 110.7 513.2 6 2015 540 110.7 513.2 24.	1	Year	2011	2012	2013	2014	2015	2	2011	678	128.8	647.5	10.5
3 Durum 128.8 126.3 117.1 124.8 110.7 4 Barley 647.5 625.2 463.7 574.7 513.2 5 Triticale 10.5 14.1 16.9 15 24.5 6 2015 540 110.7 513.2 24.	2	Wheat	678	663.3	548.4	596.1	540		_				
4 Barley 647.5 625.2 463.7 574.7 513.2 5 Triticale 10.5 14.1 16.9 15 24.5 6 2015 540 110.7 513.2 24.	3	Durum	128.8	126.3	117.1	124.8	110.7						
5 Triticale 10.5 14.1 16.9 15 24.5 6 2014 596.1 124.8 574.7 1 6 2015 540 110.7 513.2 24.								4	2013	548.4	117.1	463.7	16.9
6 2015 540 110.7 513.2 24.		-						5	2014	596.1	124.8	574.7	15
6		Triticale	10.5	14.1	16.9	15	24.5	6	2015	540	110.7	513.2	24.5
	6							7	2020			02012	2

Formatting Features

The special formatting features (merged cells, boarders, colors, bold, etc.) must be avoid as much as possible.

"Consider restructuring your data in such a way that you will not need to merge cells or other esthetic features to organize your data" (Bahlai, 2017)



Column Headers

Do not capture documentation and text descriptions in the tables themselves. The descriptive information can be recorded in the data dictionary or put into a "Note" column created for this purpose (USDA, 2017).

Ensure column headings do not contain spaces, hyphens or any other symbols. Only the underscore is allowed.

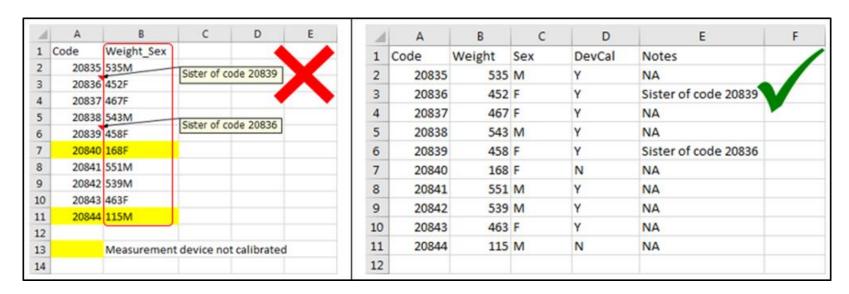
4	A	B		С	D			Α	В	С	D
2	Updates: 1	airy Production 5/03/2016		Symbol	and Space		1	Year	TempMax	N_Cattle	Milk_QTY
3							2	2000	36	766	0.8
4	Year	Maximum Temperat		N° of cattle	Quantity of Milk	_	3	2001	35	763	0.8
5	2000		36 35	766 763		0.8	4	2002	37	753	/0.8
7	2001		37	753		0.8	5	2003	38	679	
8	2003		38	679		0.7	6	2004		657	0.7
9	2004		36	657		0.7	-				
10	2005		38	685	(0.7	7	2005	38	685	0.7

[&]quot;Underscores (_) are a good alternative to spaces. Consider writing names in camel case (e.g. TestName) to improve readability" (Bahlai, 2017).

Data Entry

- 1. Data must be entered in a consistent way using always the same code for the same value.
- 2. No more than one piece of information can be in one cell.
- 3. Highlighted cells and comments must be removed since they may create problems for machine readability.

 These observations can be entered in new columns (e.g. Notes, etc.)



Note: When writing text in the cells, they can only contain text and spaces.

Null Values

Null values must be represented differently from "0". In fact, "0" corresponds to a measured data while null value means that the data is not measured at all.

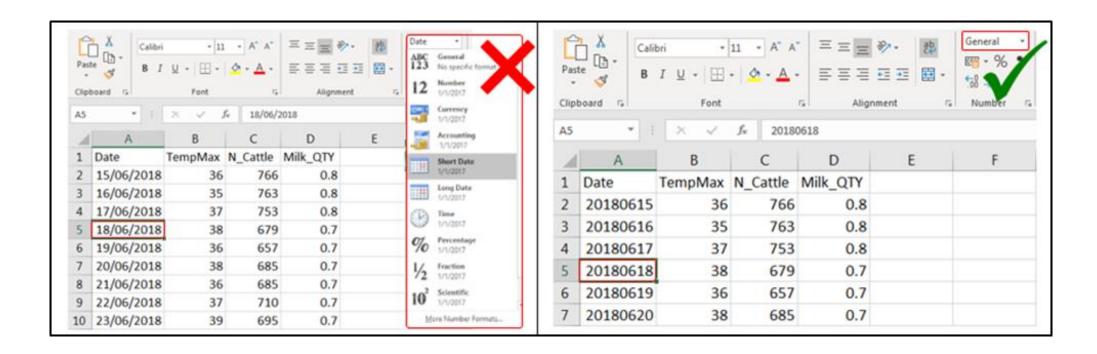
It is common to represent null values with blank cells but even to enter "NA" is a good option.

4	Α	В	С	D	E		A	Α	В	С	D	E
1	Time	Rainfall	Outflow	Total_Sed	0_200)	1	Time	Rainfall	Outflow	Total_Sed	0_200
2	0	0	0	0		0	2	0	0	0	0	0
3	10	0.8	NA	Null	None		3	10	0.8	NA	NA	NA
4	20	12.1	0	0		0	4	20	12.1	0	0	C
5	30	15.4	1.5673	0.000115	0.00000	99	5	30	15.4	1.5673	0.000115	0.0000099
6	40	8.9	9.91394	0.0003247	-		6	40	8.9	9.91394	0.0003247	NA
7	50	0	0	0		0	7	50	0	0	0	0
8	60	No Data	Missing	N/A	na		8	60	NA	NA	NA	NA
9	70	5.2	0	0		0	9	70	5.2	0	0	0
10	80	9.8	0	0		0	10	80	9.8	0	0	0

Note: As long as the missing value representation is consistent and documented, the next users can replace the choice for a null value independently (Zwicker, 2016).

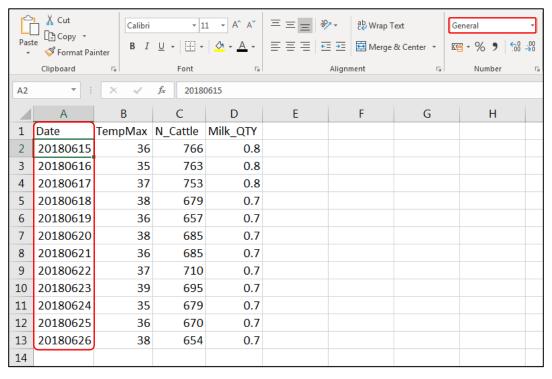
Dates and Time

The special functionalities available must not be used since they are usually guaranteed to be compatible only within the same family of products (Bahlai, 2017).



Dates and Time

Basing on ISO standards (ISO 8601:2004), the suggested format to store dates is YYYYMMDD while for time is hhmmss using the 24-hours notation (which become YYYYMMDDhhmmss when represented together).



For example, June 15, 2018 17:25:35 become 20180615172535.

Latitudes and Longitudes

The recommended standard to report latitudes and longitudes data is using decimal degrees (DD), since they guarantee the possibility of treating latitude and longitude as a simple and numeric value facilitating any next software interpretations (Callahan, 2009).

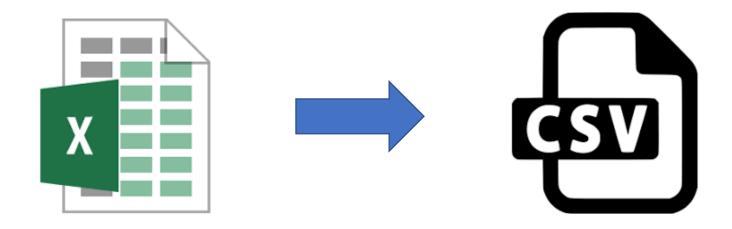
Unit	Latitude	Longitude
Decimal Degrees (DD)	40.75889	-73.98513
Degrees Minutes and Seconds (DMS)	40° 45' 32.004" N	73° 59' 6.468" W
Degrees Decimal Minutes (DM)	40° 45.5334'	-73° 59.1078′

Δ	Α	В	С
1	Site	Latitude	Longitude
2	PT-12A	39.91382	116.36363
3	PT-34B	40.75889	-73.98513
4	PT-41C	-22.90278	-43.20750
5	PT-56D	-33.86785	151.20732

Stable File Format

From Excel to CSV

It is important to save the dataset in a consistent format that can be read well into the future and is independent of changes in applications. The CSV or comma separated value format is the preferred format for most of data repositories and are the recommended one for publishing machine-readable tabular data.



Data Dictionary

Data dictionaries are used to provide detailed information about the contents of a dataset or database, such as the names of measured variables, their data types or formats, and text descriptions. In particular (Briny, 2015):

- It helps the dataset author to remember all the details about the data along the years;
- It facilitates the dataset sharing with collaborators helping them to understand and use the data files.
- It helps for personnel who are "totally unfamiliar with the data, to pick up that data, understand and reproduce the results or reuse these for new research" (Briny, 2015) activities improving the credit of the dataset.

Data Dictionary

When the data are managed in professional databases, it is possible to automatically generate data dictionaries by the available tools in the software.

While, when data are managed in spreadsheets, text files, or comma separated values, the data dictionary must be created manually. A common approach when doing it manually, is to create three main files, to save in CSV format, which contains three different level of dataset information:

- **A. Dataset Introduction:** where introductory and background explanatory information are reported;
- **B.** Dataset Elements Descriptions: where the datasets fields are listed with their related information;
- **C. Unique Identifier:** where the dataset elements and concepts are identified by dereferenceable URIs to online multilingual thesaurus.

Dataset Introduction

The purpose of the Dataset Introduction tab is to explain the contents of the dataset. Here are available the general information to make clear all the dataset aspects for next uses.

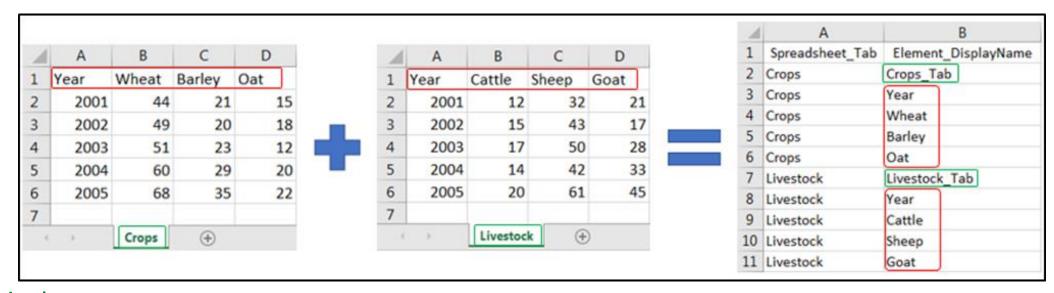
4	4	А	В	С	D	Е	F	G	Н
1	L	Description	Summary	Start_Date	End_Date	Latitude	Longitude	Author	CoAuthor
		A rich and full dataset description that	"A shorter	The date in	The date in	Latitude site	Longitude site	Dataset first	Dataset
		explains how and why it was generated	description of the	which the data	which the data	coordinate, in	coordinate, in	author.	co-author(s).
		and informs how it should be used.	dataset, usually no	collection starts	collection ends	decimal	decimal		
		Make sure that in this description are	more than a	(YYYYMMDD).	(YYYYMMDD).	degrees (DD),	degrees (DD),		
		present the experiment settings	sentence or two"			using WGS84	using WGS84		
		(location, climatic conditions, etc.), data	(USDA, 2016).			datum.	datum.		
		collection and processes methods,							
		equipment used, possible resources and							
2	2	any limiting factors (USDA, 2016).							

Dataset Elements Descriptions

This is the core of the data dictionary, since it is the document that allow the dataset users to fully understand the contents of the dataset (ORNL DAAC, 2018).

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

- **1. Spreadsheet_Tab:** The tab where is the element.
- **2. Element_DisplayName:** The dataset element name.



Dataset Elements Descriptions

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

В	C
Element_DisplayName	Description
number	Invoice number
date	Invoice date
status	Invoice status
amount	Invoice amount
customer_no	Customer number

В	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.

Source: Kononow, 2017

Note: It is important that descriptions are meaningful avoiding the text holding zero information.

Dataset Elements Descriptions

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.

\square	Α	В	С
1	Years	NetPrimaryFemale	NetPrimaryMale
2	2000	94.15018	96.6708
3	2001	96.45512	98.07967
4	2002	NA	NA
5	2003	NA	NA
6	2004	NA	NA
7	2005	98.6595	99.06998
8	2006	98.66384	98.59785
9	2007	97.41568	97.87909
10	2008	97.60895	98.38915
11	2009	98.26036	98.9663
12	2010	NA	NA
13	2011	NA	NA
14	2012	NA	NA
15	2013	NA	NA
16	2014	NA NA	NA
	← →	WorldBank_Ed	ducation (+)

4	Α	В	С
1	Years	BasicSecondary_Male	BasicSecondary_Female
2	2000	469202	493783
3	2001	497945	529867
4	2002	507290	549943
5	2003	511999	564239
6	2004	512001	572877
7	2005	505330	570187
8	2006	511128	577688
9	2007	500517	569068
10	2008	467328	538815
11	2009	447369	520339
12	2010	433814	502584
13	2011	428109	494349
14	2012	418498	490102
15	2013	408292	479153
16	2014	402896	473815
	<>	INS_StudentsEd	ucation (+)

A	Α	В	C	D	Е	F	G	Н	I
1	Spreadsheet_Tab	Element_DisplayName	Description	Units	Data_Type	Character_Length	Acceptable_Values	Required	Accepts_NullValue
2	WorldBank_Education	World Bank_Education_tab	Data about education participation in Tunisia from 1999 to 2014. Source: World Development Indicators, THE WORLD BANK. Last update 1/2/2017. Retrieved from: http://data.worldbank.org/data-catalog/world-development-indicators.	NA	NA	NA	NA	NA	NA
3	WorldBank_Education	Years	The year to which this analysis refers.	уууу	date	4	[2000,2015]	у	n
4	WorldBank_Education	NetPrimaryFemale	The element full name is "Adjusted net enrollment rate, primary, female (% of primary school age children)". Adjusted net enrollment is the number of pupils of the school-age group for primary education, enrolled either in primary or secondary education, expressed as a percentage of the total population in that age group.	%	decimals	255	NA	n	у
5	WorldBank_Education	NetPrimaryMale	The element full name is "Adjusted net enrollment rate, primary, male (% of primary school age children)". Adjusted net enrollment is the number of pupils of the school-age group for primary education, enrolled either in primary or secondary education, expressed as a percentage of the total population in that age group.	%	decimals	255	NA	n	У
6	INS_StudentsEducation	INS_StudentsEducation_tab	Data about male and female basic and secondary education in Tunisia. Last update: 17/03/2016. Source: Ministry of education (Statistique Tunisia). Retrieved from: http://www.ins.tn/en/themes/education.	NA	NA	NA	NA	NA	NA
7	INS_StudentsEducation	Years	The year to which this analysis refers.	уууу	date	4	[1990,2014]	у	n
8	INS_StudentsEducation	BasicAndSecondaryEdu_Male	The element full name is: "Number of male students in the second cycle of basic education and secondary public education". It corresponds to the male students registered in Tunisia for the different years	Individuals	numeric	6	[45000,55000]	n	У
9	INS_StudentsEducation	BasicAndSecondaryEdu_Female	The element full name is "Number of female students in the second cycle of basic education and secondary public education". It corresponds to the female students registered in Tunisia for the different years	Individuals	numeric	6	[45000,60000]	n	у

24

Source: Khawam, 2017a

Unique Identifier

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 online multilingual thesaurus.

1	Α	В	C	D	E
1	Spreadsheet tab	Element_DisplayName	Unique identifier	Source	/
2	INS_HarvestedArea	Grain	http://aims.fao.org/aos/agrovoc/c_3346	AGROVOC	
3	INS_HarvestedArea	Dried legumes	http://aims.fao.org/aos/agrovoc/c_4255	AGROVOC	
4	INS_HarvestedArea	Beans	http://aims.fao.org/aos/agrovoc/c_331566	AGROVOC	
5	INS_HarvestedArea	Root crops	http://aims.fao.org/aos/agrovoc/c_6641	AGROVOC	
6	INS_HarvestedArea	Nuts	http://aims.fao.org/aos/agrovoc/c_12873	AGROVOC	
7	INS_HarvestedArea	Fresh vegetables	http://aims.fao.org/aos/agrovoc/c_8174	AGROVOC	
8	INS_HarvestedArea	Fruits	http://aims.fao.org/aos/agrovoc/c_3131	AGROVOC	
9	INS_HarvestedArea	Citrus	http://aims.fao.org/aos/agrovoc/c_1637	AGROVOC	
10	INS_HarvestedArea	Grapes	http://aims.fao.org/aos/agrovoc/c_3359	AGROVOC	
11	INS_HarvestedArea	Olives	http://aims.fao.org/aos/agrovoc/c_12926	AGROVOC	
12	INS_HarvestedArea	Dates	http://aims.fao.org/aos/agrovoc/c_25475	AGROVOC	

Source: Khawam, 2017b

Conclusions

Dataset curation practices can be challenging depending on the status of the data. In general, they are based on a standardization of the dataset contents and the creation of the necessary documentation.

Following the practices mentioned before, starting from a dataset file in Excel format, we end with multiple files in CSV format. Thus, the result of this work is a .zip file containing the data dictionaries files and the dataset ones in CSV format.

Recommendations

All these aspects need to be considered from the first data collection steps; so that these practices become an integral part of the author's work thus reducing the heaviness of this activity as well as ensuring better results.



References

Bahlai, C., & Teal, T., (Eds). (2017, April). Data Carpentry: Data Organization in Spreadsheets Ecology lesson (Version 2017.04.0). Retrieved from http://www.datacarpentry.org/spreadsheet-ecology-lesson/

Briny, K. [University of Wisconsin Data Services]. (2015, January 23). *Data Management: Data Dictionaries*. Retrieved from https://www.youtube.com/watch?v=Fe3i9qyqPjo

Callahan, J. (2009). Standard Latitudes and Longitudes. Retrieved from http://mazamascience.com/WorkingWithData/?p=103

Khawam, H., & Najjar, D. (2017a). *Statistics on Gender and Education in Tunisia*. Retrieved from https://hdl.handle.net/20.500.11766.1/7MMLXI

Khawam, H., & Najjar, D. (2017b). Statistics on Crops, with a Focus on Barley Production and Trade, in Tunisia. Retrieved from https://hdl.handle.net/20.500.11766.1/YSTDVL

Kononow, P. (2017, August 29). *Captain Obvious' Guide to Column Descriptions - Data Dictionary Best Practices* [Blog post]. Retrieved from https://dataedo.com/blog/captain-obivous-guide-to-column-descriptions-data-dictionary-best-practices

Munoz, T., Flanders, J., Senseney, M., Davis, R., Hsu, P.H., Little, J., Jackson, L.S., Renear, A., & Trainor, K. (2017). *Frequently asked questions about data curation*. Retrieved December 5, 2018, from https://guide.dhcuration.org/faq/

Oak Ridge National Laboratory Distributed Active Archive Center (ORNL DAAC). (2018). *Data Management*. Retrieved December 5, 2018, from https://daac.ornl.gov/datamanagement/

References

Ruggles, S. (2018). *The Importance of Data Curation*. In Vannette, D., Krosnick, J. (Eds). The Palgrave Handbook of Survey Research (pp. 303-308). Palgrave Macmillan: Cham. Retrieved from https://doi.org/10.1007/978-3-319-54395-6_39

United States Department of Agriculture (USDA) (2016). *Ag Data Commons Data Submission Manual v1.3*. National Agricultural Library. Retrieved from https://data.nal.usda.gov/book/export/html/2769

United States Department of Agriculture (USDA). [National Agricultural Library]. (2017, August 9). ADC 18 - Convert data files to CSV format. Retrieved from

https://www.youtube.com/watch?v=szDWlvQOa_g&index=19&list=PL_8uALA03ZsWQ44QNKo4_ZSYSQP7gJ9h7

Zwicker, S., in Hsu, L. (2016, December 7). *How "clean" should an Excel file be to be considered machine readable*. Retrieved from https://my.usgs.gov/confluence/pages/viewpage.action?pageId=559852026

Thanks for your attention!