

# Data Organization in Spreadsheets

## Morning Session I

Jacob F. Koehler, PhD. <sup>1</sup>

<sup>1</sup>Data Carpentry

Rutgers Newark, February 11, 2018

# Keeping Track

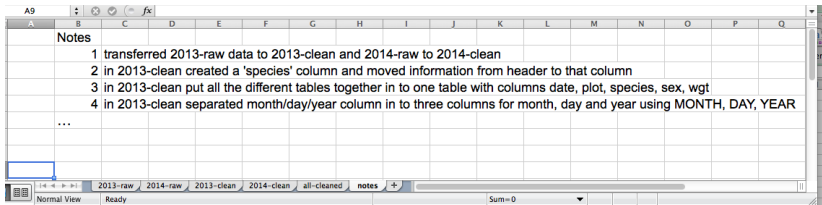


Figure: Raw Copy and Notes

## Structuring Data

Date collected	Plot	Species-Sex	Weight
1/9/78	1	DM-M	40
1/9/78	1	DM-F	36
1/9/78	1	DS-F	135
1/20/78	1	DM-F	39
1/20/78	2	DM-M	43
1/20/78	2	DS-F	144
3/13/78	2	DM-F	51
3/13/78	2	DM-F	44
3/13/78	2	DS-F	146

Figure: What Could be Better?

## Columns and Row Organization

Date collected	Plot	Species	Sex	Weight
1/9/78	1	DM	M	40
1/9/78	1	DM	F	36
1/9/78	1	DS	F	135
1/20/78	1	DM	F	39
1/20/78	2	DM	M	43
1/20/78	2	DS	F	144
3/13/78	2	DM	F	51
3/13/78	2	DM	F	44
3/13/78	2	DS	F	146

Figure: Each Column is Variable, Each Row an Observation

## Exercise I

- 1 Download the data by clicking link in lesson or in etherpad to get it from FigShare.
- 2 Open up the data in a spreadsheet program.
- 3 You can see that there are two tabs. Two field assistants conducted the surveys, one in 2013 and one in 2014, and they both kept track of the data in their own way. Now you're the person in charge of this project and you want to be able to start analyzing the data.
- 4 With the person next to you, identify what is wrong with this spreadsheet. Also discuss the steps you would need to take to clean up the 2013 and 2014 tabs, and to put them all together in one spreadsheet.

# Multiple Tables

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG		
1	lake site May 29 2012						29-May		lake site Jun 12 2012						12-Jun		lake site Jun 19 2012						19-Jun		lake site Jun 26 2012						26-Jun				
							avr	SEM																											
3	1	T1	1	1	2		T1	2.6	0.51	1	T1	6	85	99	T1	30.4	15.47126	1	T1	17	80	97		avr	SEM	1	T1	32	191	243		avr	SEM		
4	2	T1	1	2	3		T2	0.2	0.2	2	T1	9	13	21	T1	0.2	0.2	2	T1	44	126	180	T1	77.8	30.384865	2	T1	50	270	320	T1	541.6	60.913		
5	3	T1	1	3	4		control	0.2	0.2	3	T1	11	0	11		control	0.6	0.6	3	T1	18	0	18	T2	1.8	1.5620499	3	T1	6	0	6	T2	0.2	0.2	
6	4	T1	1	0	1					4	T1	0	6	6				4	T1	0	14	14		control	0.4	0.244949	4	T1	0	39	39		control	0	0
7	5	T1	0	3	3					5	T1	3	20	25				5	T1	10	70	80				5	T1	4	96	100					
8	6	T2	1	0	1					6	T2	0	0	0				6	T2	1	7	8				6	T2	0	1	1					
10	7	T2	0	0	0					7	T2	0	0	0				7	T2	0	1	1				7	T2	0	0	0					
11	8	T2	0	0	0					8	T2	1	0	1				8	T2	0	0	0				8	T2	0	0	0					
12	9	T2	0	0	0					9	T2	0	0	0				9	T2	0	0	0				9	T2	0	0	0					
13	10	T2	0	0	0					10	T2	0	0	0				10	T2	0	0	0				10	T2	0	0	0					
14	11	control	0	0	0					11	control	0	0	0				11	control	0	0	0				11	control	0	0	0					
15	12	control	0	0	0					12	control	0	0	0				12	control	0	0	0				12	control	0	0	0					
16	13	control	0	0	0					13	control	0	0	0				13	control	0	0	0				13	control	0	0	0					
17	14	control	0	0	0					14	control	0	0	0				14	control	0	1	1				14	control	0	0	0					
19	15	control	1	0	1					15	control	0	0	0				15	control	0	1	1				15	control	0	0	0					
20																																			
21	Barn site May 29 2012						29-May		Barn site Jun 12 2012						12-Jun		Barn site Jun 19 2012						19-Jun		Barn site Jun 26 2012						26-Jun				
							avr	SEM																											
22																																			
23	1	T1	3	3	6					1	T1	21	0	21				1	T1	5	0	5				1	T1	0	0	0					
24	2	T1	1	4	5		avr	SEM	2	T1	36	74	110		avr	SEM	2	T1	65	502	567		avr	SEM	2	T1	44	2057	2101	T1	431.8	417.33			
25	3	T1	0	0	0		T1	2.4	1.288	3	T1	13	0	13	T1	50.6	20.10124	3	T1	10	7	17	T1	119.4	111.92882	3	T1	12	20	32	T2	0.4	0.4		
26	4	T1	0	0	0		T2	0.4	0.245	4	T1	7	0	7	T2	1	0.774597	4	T1	0	6	6	T2	5	2.1908902	4	T1	0	16	16		control	1.2	0.5831	
27	5	T1	0	1	1		control	1	0.316	5	T1	1	0	2		control	2.2	1.714643	5	T1	0	2	2		control	2.8	0.969556	5	T1	0	10	10			
28	6	T2	0	0	0					6	T2	1	0	1				6	T2	0	8	8				6	T2	0	0	0					
29	7	T2	0	0	0					7	T2	0	4	4				7	T2	0	12	12				7	T2	0	0	0					
30	8	T2	0	1	1					8	T2	0	0	0				8	T2	0	0	0				8	T2	0	0	0					
31	9	T2	0	1	1					9	T2	0	0	0				9	T2	0	0	0				9	T2	0	0	0					
32	10	T2	0	0	0					10	T2	0	0	0				10	T2	0	0	0				10	T2	0	0	0					
33	11	control	0	0	0					11	control	0	1	1				11	control	0	5	5				11	control	0	2	2					
34	12	control	1	1	1					12	control	1	1	2				12	control	1	1	2				12	control	1	0	1					
35	13	control	1	1	1					13	control	0	0	0				13	control	0	0	0				13	control	0	0	0					
36	14	control	1	1	1					14	control	1	9	9				14	control	0	5	5				14	control	0	1	1					
37	15	control	2	2	2					15	control	0	1	1				15	control	0	2	2				15	control	1	0	0					
38																																			
39																																			

Figure: Recall Column Organization

# Using Problematic Null Values

**Table 1.** Commonly used null values, limitations, compatibility with common software and a recommendation regarding whether or not it is a good option. Null values are indicated as compatible with specific software if they work consistently and correctly with that software. For example, the null value "NULL" works correctly for certain applications in R, but does not work in others, so it is not presented in the table as R compatible.

Null values	Problems	Compatibility	Recommendation
0	Indistinguishable from a true zero		Never use
Blank	Hard to distinguish values that are missing from those overlooked on entry. Hard to distinguish blanks from spaces, which behave differently.	R, Python, SQL	Best option
-999, 999	Not recognized as null by many programs without user input. Can be inadvertently entered into calculations.		Avoid
NA, na	Can also be an abbreviation (e.g., North America), can cause problems with data type (turn a numerical column into a text column). NA is more commonly recognized than na.	R	Good option
N/A	An alternate form of NA, but often not compatible with software		Avoid
NULL	Can cause problems with data type	SQL	Good option
None	Uncommon. Can cause problems with data type	Python	Avoid
No data	Uncommon. Can cause problems with data type, contains a space		Avoid
Missing	Uncommon. Can cause problems with data type		Avoid
~,+,,	Uncommon. Can cause problems with data type		Avoid

Figure: Null Suggestions

## Using Formatting to Convey Information

Plot: 2			
Date collect	Species	Sex	Weight
1/8/14	NA		
1/8/14	DM	M	44
1/8/14	DM	M	38
1/8/14	OL		
1/8/14	PE	M	22
1/8/14	DM	M	38
1/8/14	DM	M	48
1/8/14	DM	M	43
1/8/14	DM	F	35
1/8/14	DM	M	43
1/8/14	DM	F	37
1/8/14	PF	F	7
1/8/14	DM	M	45
1/8/14	OT		
1/8/14	DS	M	157
1/8/14	OX		
2/18/14	NA	M	218
2/18/14	PF	F	7
2/18/14	DM	M	52
	measurement device not calibrated		

Figure: Highlighting Problematic Data



## Using Formatting to Convey Information

Plot: 2				
Date collect	Species	Sex	Weight	
1/8/14	NA			
1/8/14	DM	M	44	
1/8/14	DM	M	38	
1/8/14	OL			
1/8/14	PE	M	22	
1/8/14	DM	M	38	
1/8/14	DM	M	48	
1/8/14	DM	M	43	
1/8/14	DM	F	35	
1/8/14	DM	M	43	
1/8/14	DM	F	37	
1/8/14	PF	F	7	
1/8/14	DM	M	45	
1/8/14	OT			
1/8/14	DS	M	157	
1/8/14	OX			
2/18/14	NA	M	218	
2/18/14	PF	F	7	
2/18/14	DM	M	52	
	measurement device not calibrated			

Date collect	Species	Sex	Weight	Calibrated
1/8/14	NA			
1/8/14	DM	M	44	Y
1/8/14	DM	M	38	Y
1/8/14	OL			
1/8/14	PE	M	22	Y
1/8/14	DM	M	38	Y
1/8/14	DM	M	48	Y
1/8/14	DM	M	43	Y
1/8/14	DM	F	35	Y
1/8/14	DM	M	43	Y
1/8/14	DM	F	37	Y
1/8/14	PF	F	7	Y
1/8/14	DM	M	45	Y
1/8/14	OT			
1/8/14	DS	M	157	N
1/8/14	OX			
2/18/14	NA	M	218	N
2/18/14	PF	F	7	Y
2/18/14	DM	M	52	Y

Figure: Highlighting Problematic Data

Figure: Solution: New Column

# Problematic Field Names

Good Name	Good Alternative	Avoid
Max_temp_C	MaxTemp	Maximum Temp (°C)
Precipitation_mm	Precipitation	precmm
Mean_year_growth	MeanYearGrowth	Mean growth/year
sex	sex	M/F
weight	weight	w.
cell_type	CellType	Cell Type
Observation_01	first_observation	1st Obs

Figure: Example Names



## Dates as Data

**Challenge:** pulling month, day and year out of dates

- In the **dates** tab of your spreadsheet you have the data from 2014 plot 3. There's a **Date collected** column.
- Let's extract month, day and year from the dates to new columns. For this we can use the built in Excel functions **YEAR()** **MONTH()** **DAY()**

## Hours, Minutes, and Seconds

Current time and date are best retrieved using the functions **NOW()**, which returns the current date and time, and **TODAY()**, which returns the current date. The results will be formatted according to your computer's settings.

- 1 Extract the year, month and day from the current date and time string returned by the **NOW()** function.
- 2 Calculate the current time using **NOW()-TODAY()**.
- 3 Extract the hour, minute and second from the current time using functions **HOUR()**, **MINUTE()** and **SECOND()**.
- 4 Press **F9** to force the spreadsheet to recalculate the **NOW()** function, and check that it has been updated.

## Preferred date format

	A	B	C	D	E	F	G	H	I
1	What I typed in	day-month	DOW, month, day, year	month-year	Initial-year	M/D/YYYY	DD/MM/YYYY	DD/MM/YY	number
2	2-Jul	2-Jul	Wednesday, July 02, 2014	Jul-14	J-14	7/2/2014	02/07/2014	07/02/14	41822
3	Jul-14	14-Jul	Monday, July 14, 2014	Jul-14	J-14	7/14/2014	14/07/2014	07/14/14	41834
4	1-Jan-1900	1-Jan	Sunday, January 01, 1900	Jan-00	J-00	1/1/1900	01/01/1900	01/01/00	1

Figure: Year, Month, Day in Separate Columns

## Preferred date format

	A	B	C	D	E	F	G	H	I
1	What I typed in	day-month	DOW, month, day, year	month-year	Initial-year	M/D/YYYY	DD/MM/YYYY	DD/MM/YY	number
2	2-Jul	2-Jul	Wednesday, July 02, 2014	Jul-14	J-14	7/2/2014	02/07/2014	07/02/14	41822
3	Jul-14	14-Jul	Monday, July 14, 2014	Jul-14	J-14	7/14/2014	14/07/2014	07/14/14	41834
4	1-Jan-1900	1-Jan	Sunday, January 01, 1900	Jan-00	J-00	1/1/1900	01/01/1900	01/01/00	1

Figure: Year, Month, Day in Separate Columns

**Challenge:** What happens to the dates in the “dates” tab of our workbook if we save this sheet in Excel (in csv format) and then open the file in a plain text editor (like TextEdit or Notepad)? What happens to the dates if we then open the csv file in Excel?

# Quality Assurance

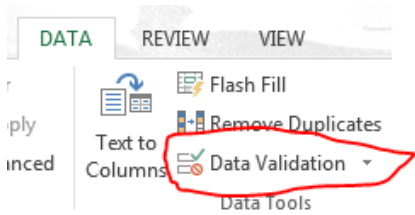


Figure: Checking Values are Valid During Data Entry



# Quality Control

## Sorting

We've combined all of the tables from the messy data into a single table in a single tab. Download this semi-cleaned data file to your computer.

Once downloaded, sort the Weight\_grams column in your spreadsheet program from Largest to Smallest.

What do you notice?

# Conditional Formatting

## Challenge:

- 1 In the main Excel menu bar, click **Format > Conditional Formatting...** Click the + to add a formatting rule.
- 2 Apply a **2-Color Scale** formatting rule with the lowest values set to orange and the highest values set to yellow.
- 3 Now we can scan through and different colors will stand out. Do you notice any strange values?

# Exporting Data

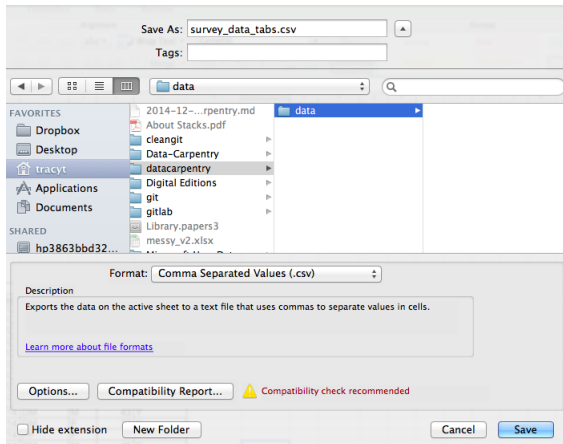


Figure: Saving as .csv