

Collaboration in the 21st century: tips and tools for dealing with data

(or, how I learned to stop worrying and love version control)

Nick Waters

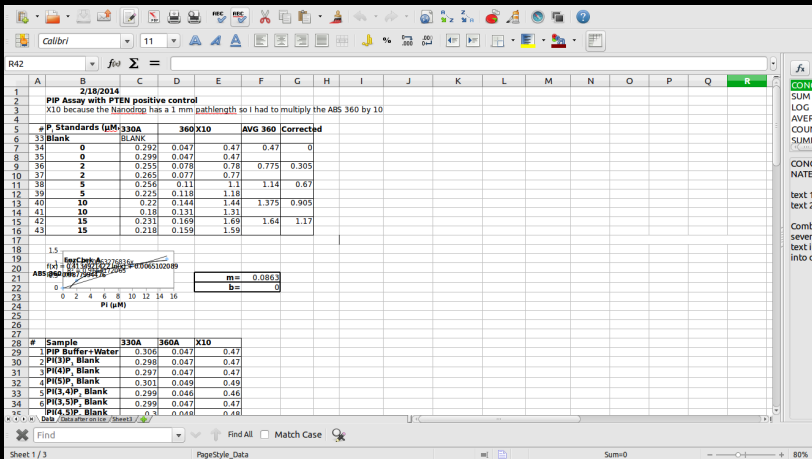
Hot Topic: March 7, 2017

National University of Ireland, Galway
James Hutton Institute, Dundee

The Problem: other people's data

The Problem: my data

Bad Data



Bad Data



Q75 Colibri 11 \sum = S3.2

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
1														
2														
3							X Xnum >	###	1/1/1900	1/2/1900	4	5	6	7
4							a 1 1 10	53.7	55.9	53.3	55	53	54.6	53.2
5							b 2 1 10	51.8	51.3	53.7	51.5	49.4	50.3	50.2
6							c 3 1 10	53.2	53.4	62	50.1	60.8	55.5	54.3
7							d 4 1 10	58.4	52	48.7	53.6	60.2	57.3	55.6
8							e 5 1 10	53.1	52.6	49.7	54.7	61.4	54.3	52.7
9							f 6 1 10	54	49.2	52.5	58.1	63.2	55.3	50.4
10							g 7 1 10	56.6	52.9	55.3	50.6	53.3	55.4	53.4
11							h 8 1 10	62.4	55	52.6	52.4	65.6	61	55.6
12							i 9 1 10	48.7	54.4	57.4	47.8	57.7	55.5	49.7
13	Instruction						j 10 1 1	55.1	52.9	55.8	48.8	50.1	55.9	49.6
14	1 copy table into this document						k 11 1 1	57.8	62	55.2	53.9	57.2	55.4	60.1
15	2 Concatenate row descriptions into single row with a						l 12 1 1	61.2	59.1	48.2	52.2	50.6	56.5	53.4
16	3 click button, follow instructions						m 13 1 1	51.3	61	57.6	51.9	49	55.3	55.8
17	4						n 14 1 1	54	51	56	56.3	50	54.3	53.4
18	5						o 15 1 1	55.4	51.8	60.9	62.8	51.1	56.6	54.5
19							p 16 1 1	48.3	50.9	57.3	51.3	59.1	50.4	54.9
20							q 17 1 1	52	58.7	49.9	54.6	58.2	52	54.2
21							r 18 1 1	62.4	52.3	63	48.7	55.4	54.3	56.2
22							s 19 1 1	51.7	51.9	63.9	60.2	57.2	54.2	56.3
23							t 20 1 1	49.4	47.6	53.8	58.9	51.3	52	49.9
24							u 21 1 1	65	54.3	55.8	53.5	55.2	56.3	54.3
25							v 22 1 1	56.6	51.8	52.4	51	61	50.3	52.3
26							w 23 1 1	51.1	67	60.3	54.2	60.7	52.3	52.1
27							x 24 1 1	50.5	53.5	53.9	54	50.9	52.3	51.6
28							a 1 3 30	53.5	54.3	54.2	57.1	49.6	56.2	54.1
29							b 2 3 30	57.1	54.3	52.1	54.3	54.1	57.1	54.2
30							c 3 3 30	57.6	54.1	55.3	54.6	58.2	57	58

Sheet1 (2) / Sheet2 / Sheet3 /

Find Find All Match Case

Sheet 1 / 3 PageStyle_Sheet1 (2) Sum=S3.2 100%

Bad Data



	100		25			
	500		125			
	337	372	543	506	577	629
400 400	4.90	4.21	2.37	2.31	1.57	5.84
11/05	.05	.05	.05	.05	.05	.05
9.00	.434	.31	.59	.926	.707	.27
11.50	.42	.602	9.76	1.415	1.28	1.4
12.25	5.43	8.37	1.38	.89	.17	.89
140	.282 #1	.233	.226	.328	.244	.259
155	.409	.323	.324	.470	.317	.369
207	.506 #2	.406	.418	.552	.421	.472
220	.608 #3	.49	.514	.743	.561	.608
234	.921 #4	.77	.702	.982	.772	.874
250	.122 #5	.978	934	1396	1004	1590
205	.118 #6	.118	11.106	1550	12.46	1396

The Real Problem: Separation of Concerns



Separation of concerns: The practice of separating data from analysis.

Why do we need to separate data from analysis?

- Saves time for repeated analysis
- Reproducible results
- Prevents data loss or contamination
- Prevents user error
- Easier to track over time



Raw data should never be changed.

- Store it in plain text format in utf-8
- Store it with metadata in the same directory also in plain text
- Make sure it is backed up



Do's:

- Include a “metadata” sheet
- Make new columns/rows for each step in analysis
- Document each step
- Explicitly set data types

Don'ts:

- Link Excel workbooks or sheets
- Rely on color coding
- Use macros
- Copy and paste



Literate programming mixes text description with data analysis. Examples include Jupyter notebooks, knitr/sweave, emacs's org-mode.

Do's:

- Explicitly state where the data is coming from
- Keep the analysis scoped to a single folder
- Document each step
- Ensure compilation will fail with ill-formatted data
- When finished, zip the entire directory for stable storage.

Don'ts:

- Mix literate programming with hardcoding
- Leave out print statements
- Forget sanity checks



R, python, MATLAB, etc

Do's:

- Use version control
- Explicitly state where the data is coming from
- Write tests and comments for all the interesting steps
- Make the program write out a log
- Ensure compilation will fail with ill-formatted data
- Time-stamp results

Don'ts:

- Hard-code paths
- Write unhelpful comments
- Write similar code
- Forget sanity checks
- Rely on Stack Overflow

Version control



Version control keeps track of changes made to data and other files

- Dropbox/Google Drive
- Google docs
- Time Machine
- Lab wiki's
- Git, subversion, or mercurial

Things that should be version controlled



- Protocols
- Manuscripts
- Analysis pipeline

Things that shouldn't be version controlled



- Data: backups are different from version control
- Temporary files
- Complex file formats, big files, etc

Takeaways



Write data for computers, write code for people



- Identify analyses that could be automated
- If you use excel, make a “metadata” sheet
- Never link Excel workbooks
- Restructure project folders