

Reproducibility

Rosana Zenil-Ferguson

4/5/2021

Useful advice

- ▶ Your closest collaborator is you six months ago, and you do not reply to emails -*Karl Broman*
- ▶ The most important tool is the mindset, when starting, that the end product will be reproducible -*Keith Baggerly*
- ▶ Assume that everything that you are doing right now will need to be redone at some point in the future: be prepared

Five stages of reproducibility

(From Claudia Solis-Lemus PhD)

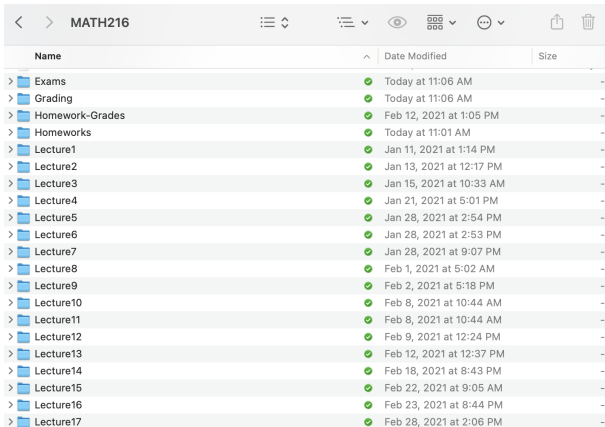
1. Denial: I do not need to be reproducible. I have not kept track of code/scripts in years and I have been just fine. People exaggerate. We do not have to be that paranoid
2. Anger: Why do I have to write these stupid notes!? It takes twice the time to write notes and do the work. I could simply do the work! This is stupid and ridiculous! I am just wasting my time with notes and comments that nobody cares about!
3. Bargaining: Well, perhaps it is ok if I only keep notes in the very final script or the very final function. That makes sense. No one needs to know or would even care to read my other code. Yes, maybe it is ok if I only comment at the end on the project

Five stages of reproducibility

(From Claudia Solis-Lemus PhD)

4. Depression: I do not understand my notes. The comments that I made a year ago do not mean anything to me anymore. This has totally failed. I am a reproducibility failure. If I am not able to understand my own notes, no one will
5. Acceptance: I understand that being reproducible is a process. No one does this right the first time. No one does it right period. We are all learning, and all I can do is try my best to make notes/comments and be honest and open about my research process

Basic organizational skills- Directories



MATH216			
Name		Date Modified	Size
> Exams	✓	Today at 11:06 AM	--
> Grading	✓	Today at 11:06 AM	--
> Homework-Grades	✓	Feb 12, 2021 at 1:05 PM	--
> Homeworks	✓	Today at 11:01 AM	--
> Lecture1	✓	Jan 11, 2021 at 1:14 PM	--
> Lecture2	✓	Jan 13, 2021 at 12:17 PM	--
> Lecture3	✓	Jan 15, 2021 at 10:33 AM	--
> Lecture4	✓	Jan 21, 2021 at 5:01 PM	--
> Lecture5	✓	Jan 28, 2021 at 2:54 PM	--
> Lecture6	✓	Jan 28, 2021 at 2:53 PM	--
> Lecture7	✓	Jan 28, 2021 at 9:07 PM	--
> Lecture8	✓	Feb 1, 2021 at 5:02 AM	--
> Lecture9	✓	Feb 2, 2021 at 5:18 PM	--
> Lecture10	✓	Feb 8, 2021 at 10:44 AM	--
> Lecture11	✓	Feb 8, 2021 at 10:44 AM	--
> Lecture12	✓	Feb 9, 2021 at 12:24 PM	--
> Lecture13	✓	Feb 12, 2021 at 12:37 PM	--
> Lecture14	✓	Feb 18, 2021 at 8:43 PM	--
> Lecture15	✓	Feb 22, 2021 at 9:05 AM	--
> Lecture16	✓	Feb 23, 2021 at 8:44 PM	--
> Lecture17	✓	Feb 28, 2021 at 2:06 PM	--

Figure 1: My math class folder

Basic organizational skills- Directories












>  images			Jan 28, 2021 at 7:08 PM
	Lecture8_post.pdf		Feb 1, 2021 at 4:59 AM
	Lecture8_pre.pdf		Feb 1, 2021 at 5:00 AM
	Lecture8.mp4		Jan 29, 2021 at 10:28 AM
	Lecture8.pdf		Jan 29, 2021 at 9:17 AM
	Lecture8.Rmd		Jan 29, 2021 at 9:17 AM

Figure 2: Lecture folder

Tracking everything you do with Rmarkdown

For example, I like it for my class and tutorials because it can write math

```
1 ---
2 title: "Lecture 4"
3 author: "MATH 216"
4 date: "1/20/2021"
5 output: beamer_presentation
6 ---
7
8 ## Derivatives
9 In MATH 215 you learned that a derivative is no other thing that the slope of a tangent
10 line. In biology, we interpret derivatives as the rate of change at a given point.
11
12 The formal definition of the derivative is
13
14 
$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}$$

15
16
17 However, for the purpose of this class we will just revise the essential derivatives and
18 rules.
19 ## Tangent of a function-Rate of change
```

Derivatives
Tangent of a func...
Interpreting deriv...
Example 1. Deriv...
Example 1.
Basic differentiati...
Basic differentiati...
Basic differentiati...
New derivatives f...
New derivatives f...
New derivatives f...
Example 4: The r...
Example 4: How ...
New derivatives f...

Figure 3: Writing math in Rmarkdown

Tracking everything you do with Rmarkdown

It can also accumulate all your code

```
65 ▾ ## Example 1. Covid-19 pandemic: Exponential Growth
66 When the speed of growth is *proportional to the size of the population*, that's
    exponential growth
67 ▾ ```{r, message=FALSE}
68 Date<-seq(as.Date("2020-02-23"), as.Date("2020-03-23"), by="days")
69 Cases<- c(10,11,13,18,22,30,42,47,69,109,164,220,271,352,412,469,617,876,1292,1766,2244,26
    05,3047,3657,4427,5426,6479,7738,8934,10312)
70 covid_UK<-data.frame(Date,Cases)
71 ▴
72
73 ▾ ```{r, echo=FALSE, results='asis'}
74 library(knitr)
75 kable(covid_UK[1:10,],caption="Covid-19 cases in the UK")
76 ▴
77
78 ▾ ```{r, message=FALSE}
79 log.covid_UK<-data.frame(Date,log(Cases))
80 ▴
81 ▾ ## Example 1. Covid-19 pandemic: Exponential Growth
```

Models of Popula...
Exponential growth
Example 1. Covid...
Example 1. Covid...
Example 2. Diver...
Logistic Growth d...
Example: Logistic...
Example: Logistic...

Figure 4: R code in Rmarkdown

To git or not to git

- ▶ Personally: I don't have a preference. I think you can be as organized as possible in folders as long as you comment your code
- ▶ If you don't use git as often, use dropbox- \$10 a month (I don't think about backups)

The basic data vs. the data transform







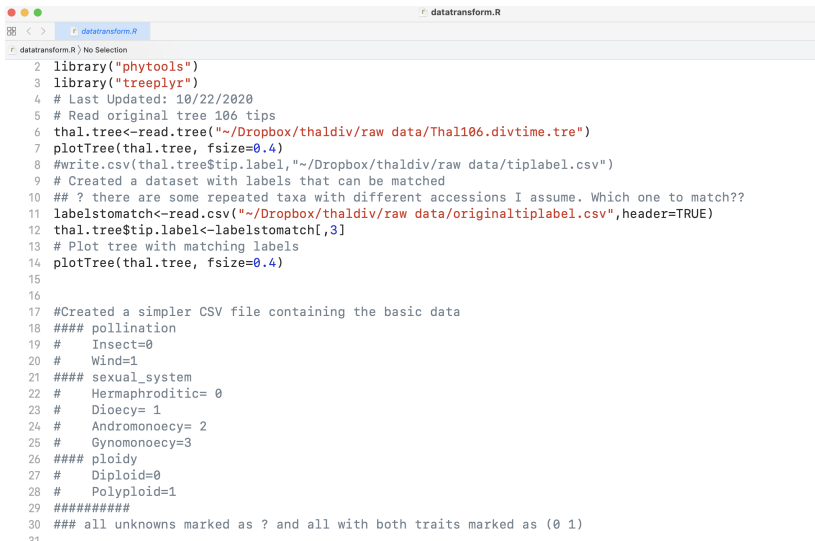
-
- >  basic data
 - >  code for Veronica
 - >  figures
 - >  raw data
 - >  results
 - >  rev code

Figure 5: Raw data folder: Do not touch!

The data transform file



```
1 library("phytools")
2 library("treeplyr")
3 # Last Updated: 10/22/2020
4 # Read original tree 106 tips
5 thal.tree<-read.tree("~/Dropbox/thaldiv/raw data/Thal106.divtime.tre")
6 plotTree(thal.tree, fsize=0.4)
7 #write.csv(thal.tree$tip.label, "~/Dropbox/thaldiv/raw data/tiplabel.csv")
8 # Created a dataset with labels that can be matched
9 ## ? there are some repeated taxa with different accessions I assume. Which one to match??
10 labelstomatch<-read.csv("~/Dropbox/thaldiv/raw data/originaltiplabel.csv",header=TRUE)
11 thal.tree$tip.label<-labelstomatch[,3]
12 # Plot tree with matching labels
13 plotTree(thal.tree, fsize=0.4)
14
15
16
17 #Created a simpler CSV file containing the basic data
18 ##### pollination
19 #   Insect=0
20 #   Wind=1
21 ##### sexual_system
22 #   Hermaphroditic= 0
23 #   Dioecy= 1
24 #   Andromonoecy= 2
25 #   Gynomoecy=3
26 ##### ploidy
27 #   Diploid=0
28 #   Polyploid=1
29 #####
30 ### all unknowns marked as ? and all with both traits marked as (0 1)
```

Figure 6: datatransform.R