

# Data Carpentry

Day 2

# Spreadsheets

- Make it a rectangle
- Rows = observations, columns = variables
- One head row; avoid spaces
- One data type per cell
- Fill in all cells
- Consistently code missing values
- Care about date data
- Don't do calculations in raw data files
- Save as CSV files
- Don't use font color or highlighting to code data

# OpenRefine

- For cleaning and exploration of data
- NOT for editing your raw data!
- Use Facets and filters to explore
- Split columns
- Remove training/ending text
- Find outliers
- All actions are reproducible

# SQL

- SELECT (choose columns)
- FROM (data sheet(s))
- WHERE (subset specific observations)
- AND/OR/IN (used in setting criteria)
- ORDER BY (sort data)
- GROUP BY (lump data into groups)
- COUNT & SUM (summarization)
- JOIN ON (combining data)

# dplyr

## R function

- select
- filter
- mutate
- group\_by
- summarize
- arrange

## SQL Keyword

SELECT

WHERE

(weight/1000)

GROUP BY

COUNT, AVG, SUM

ORDER BY

“File organization and naming are powerful weapons against chaos.”

-Jenny Bryan

# Organizing projects

- All files in common folder (directory)
- Separate raw data from “clean” data
- Separate code (and output) from data
- Use file names that are meaningful, sortable, & consistent
- Code dates: 2017-01-11
  - raw\_data/
  - in\_process\_data/
  - clean\_data/
  - code/
  - reports/

“Your closest collaborator is you from six months ago, but you don’t reply to emails.”

-(Paraphrasing) Mark Holder

Have sympathy for your future self--be an organized analyst!



# Today: R!

- Full programming language
- Focused on programming and data
- Super for data analysis and visualization
- Great community of supporters
- R Archive has >9000 add-on packages
- RStudio: “Integrated Development Environment” (IDE) for R

# Challenge

What would  $y$  equal after these three lines of code were executed (try to answer without running them first!)? Why? How would you make it equal something else?

```
 $x \leftarrow 50$ 
```

```
 $y \leftarrow x * 2$ 
```

```
 $x \leftarrow 75$ 
```

# Challenge

Use the *nrow()* function + indexing to save just the last row of *surveys* into a new object called *surveys\_last*

# Think about it

Why doesn't *mean(heights, TRUE)* work? Hint: Check *mean*'s help page!

# Challenge

With your neighbors, make a scatterplot of *weight* vs. *hindfoot\_length* for just *species\_id* “DM.”

# Challenge

1. Get the **mean** *weights* and *hindfoot\_lengths* for every species
2. Get **counts** for each species
3. Make a **scatterplot** of mean *weights* vs. mean *hindfoot\_lengths*
4. Make the point size vary by the counts.

# Challenge

Make a line graph of the *counts* over *time* of just two species: “*DM*” and “*DS*.” Add the points to the plot.

Hint: `filter(species_id == “DM “ | species_id == “DS”)`

| is OR in R

# Challenge

- Now, plot species counts over time such that there is a different facet for each species, and in each, the data are separated by sex with colors.