# A Field Guide to Base R

Josh Allen

Department of Political Science at Georgia State University

1/27/23

# A Review of the Basics of R

# Setting Your Working Directory

- Your working directory is where all your files live

- You may know where your files are…

- But R does not

- If you want to use any data that does not come with a package you are going to need to tell R where it lives

# Cats and Boxes



- You **can** put a box inside a box

- You **can** put a cat inside a box

- You **can** put a cat inside a box inside of a box

- You **cannot** put a box inside a cat

- You **cannot** put cat in a cat

# Working Directories

```
1  getwd()
```

```
[1]
"/Users/josh/Dropbox/Research-
Data-Services-Workshops/8810-
guest-lecture"
```

```
1  setwd("path/to/your/project") #mac/linux
2  setwd("path\to\your\project") # windows
```

# How To Make Your Life Easier



source: Jenny Bryan
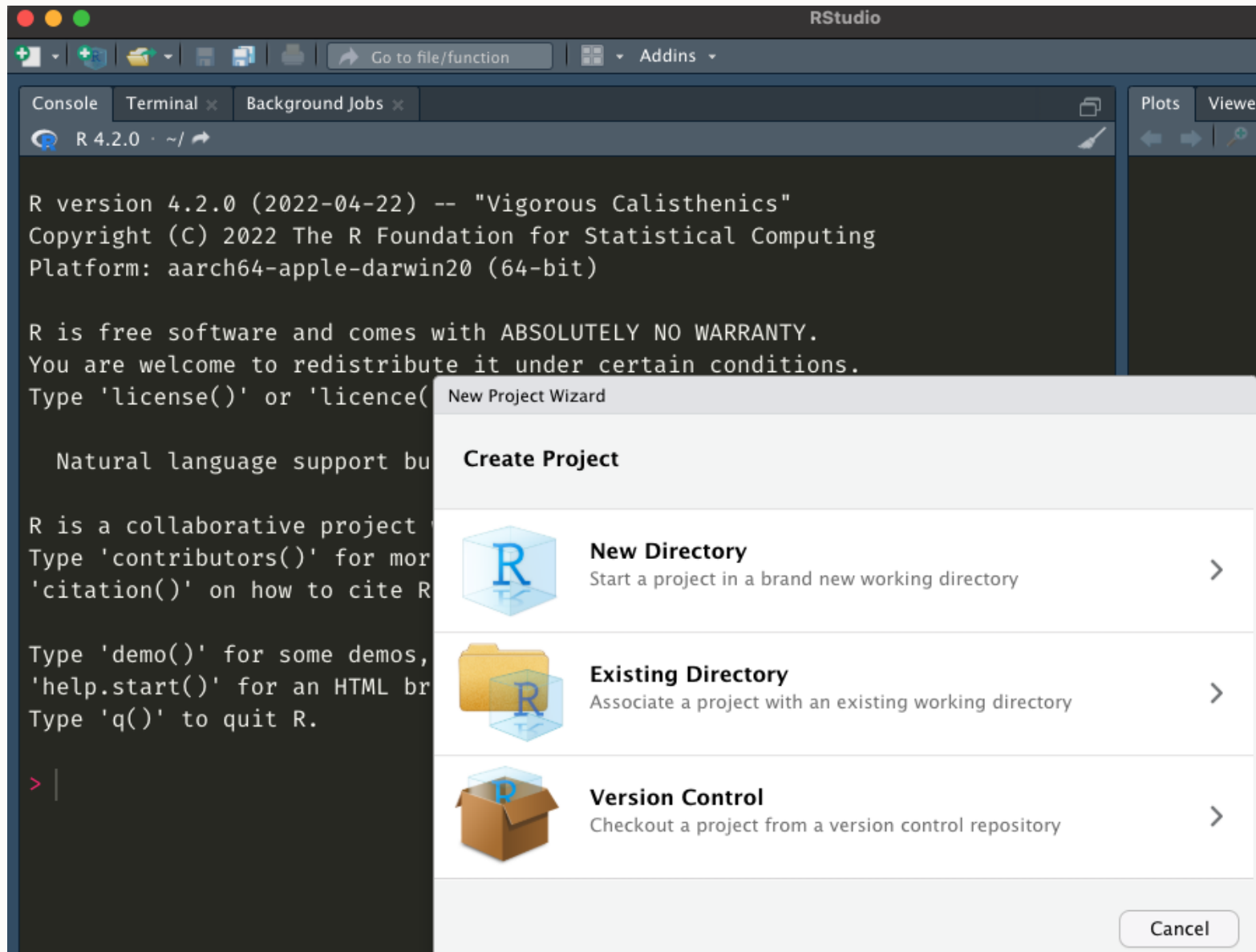
# How To Make Your Life Easier

## Working Directory for My Laptop

```
"/Users/josh/Dropbox/Research-Data-
Services-Workshops/research-data-
services-r-workshops/slides"
```

## Working Directory of My Office Computer

```
"/Volumes/6TB Raid
10/Dropbox/Research-Data-Services-
Workshops/research-data-services-r-
workshops/slides"
```

# R Projects

# The Mantra

- Everything in R is an Object

- Everything has a name

- You do stuff with functions

- Packages(i.e. libraries) are homes to pre-written functions.

    - You can also write your own functions and in some cases should.

# An Example

```r
1 digi <- c("1","2","3","4")
2 mean(digi)
```

```
[1] NA
```

```r
1 numbs <- c(1:4)
2 mean(numbs)
```

```
[1] 2.5
```

```r
1 class(digi)
```

```
[1] "character"
```

```r
1 lets <- letters
2 class(lets)
```

```
[1] "character"
```

R Some Basics

# Basic Maths

- R is equipped with lots of mathematical operations

```
1  2+2 ## addition
```

[1] 4

```
1  4-2 ## subtaction
```

[1] 2

```
1  600*100 ##multiplication
```

[1] 60000

```
1  100/10 ##division
```

[1] 10

```
1  10*10/(3^4*2)-2 ## Pemdas
```

[1] -1.382716

```
1  log(100)
```

[1] 4.60517

```
1  sqrt(100)
```

# Basic Maths

R is also equipped with modulo operations (integer division and remainders), matrix algebra, etc

```r
1  100 %/% 60 # How many whole hours in 100 minutes?
```

```
[1] 1
```

```r
1  100 %% 60 # How many minutes are left over?
```

```
[1] 40
```

```r
1  m <- matrix(1:8, nrow=2)
2  n <- matrix(8:15, nrow=4) # this is just me creating matrices
3  mat <- matrix(1:15, ncol = 5)
4  m %*% n # Matrix multiplication
```

```
      [,1] [,2]
[1,]  162  226
[2,]  200  280
```

```r
1  t(mat) # transpose a matrix
```

```
      [,1] [,2] [,3]
[1,]     1    2    3
[2,]     4    5    6
```

```
[3,]     7    8    9
[4,]    10   11   12
[5,]    13   14   15
```

# Logical Statements & Booleans

| Test | Meaning | Test | Meaning |
|------|---------|------|---------|
| x < y | Less than | x %in% y | In set |
| x > y | Greater than | is.na(x) | Is missing |
| == | Equal to | !is.na(x) | Is not missing |
| x <= y | Less than or equal to | | |
| x >= y | Greater than or equal to | | |
| x != y | Not equal to | | |
| x | y | Or | | |
| x & y | And | | |

# Booleans and Logicals in Action

```
1  1>2
```

[1] FALSE

```
1  1<2
```

[1] TRUE

```
1  1 == 2
```

[1] FALSE

```
1  1 < 2 | 3 > 4 ## only one test needs to true to return true
```

[1] TRUE

```
1  1 < 2 & 3>4 ## both tests must be true to return true
```

[1] FALSE

# Logicals, Booleans, and Precedence

- `R` like most other programming languages will evaluate our logical operators(`==`, `>`, etc) before our booleans(`|`, `&`, etc).

```
1  1 > 0.5 & 2
```
```
[1] TRUE
```

- What's happening here is that R is evaluating two separate "logical" statements:

- `1 > 0.5`, which is is obviously TRUE.

- `2`, which is TRUE(!) because R is "helpfully" converting it to `as.logical(2)`.

- It is way safer to make explicit what you are doing.

- If your code is doing something weird it might just be because of precedence issues

  - See R Cookbook 2.11

```
1  1 > 0.5 & 1 > 2
```
```
[1] FALSE
```

# Other Useful Tricks

Value matching using `%in%`

To see whether an object is contained within (i.e. matches one of) a list of items, use `%in%`.

```
1  4 %in% 1:10
```

```
[1] TRUE
```

```
1  4 %in% 5:10
```

```
[1] FALSE
```

# Cool Now What?

- While this is boring it opens up lots

- We may need to set up a group of tests to do something to data.

- We may need all this math stuff to create new variables

- However we need to *Assign them* to reuse them later in functions.

  - Including datasets

# Everything is an Object

# What are Objects?

- Objects are what we work with in R

```
 [1] "is.array"                  "is.atomic"
 [3] "is.call"                   "is.character"
 [5] "is.complex"                "is.data.frame"
 [7] "is.double"                 "is.element"
 [9] "is.environment"            "is.expression"
[11] "is.factor"                 "is.finite"
[13] "is.function"               "is.infinite"
[15] "is.integer"                "is.language"
[17] "is.list"                   "is.loaded"
[19] "is.logical"                "is.matrix"
[21] "is.na"                     "is.na.data.frame"
[23] "is.na.numeric_version"     "is.na.POSIXlt"
[25] "is.na<-"                   "is.na<-.default"
[27] "is.na<-.factor"            "is.na<-.numeric version"
```

# Vectors

- Come in **two flavors**

- **Atomic**: all the stuff must be the same type

- **Lists**: stuff can be different types

```r
1  my_vec <- c(1:10)
2  is.vector(my_vec)
```

```
[1] TRUE
```

```r
1  my_list <- list(a = c(1:4), b = "Hello World", c = data.frame(x = 1:10, y = 1:10))
2  is.vector(my_list)
```

```
[1] TRUE
```

# Atomic Vectors

- Come in a variety of flavors

- Numeric: Can contain whole numbers or decimals

- Logicals: Can only take two values TRUE or FALSE

- Factors: Can only contain predefined values. Used to store categorical data

  - Ordered factors are special kind of factor where the order of the level matters.

- Characters: Holds character strings

  - Base R will often convert characters to factors. That is bad because it will choose the levels for you

# Lists

- Lists are everywhere in R

```r
1  data_frame <- data.frame(a = rnorm(3),
2                           b = rnorm(3))
3  typeof(data_frame)
```

```
[1] "list"
```

```r
1  dats_wrong <- data.frame(a = 1:3,
2                           b = 1:4)
```

```
Error in data.frame(a = 1:3, b = 1:4): arguments imply differing
number of rows: 3, 4
```

```r
1  example_mod <- lm(body_mass_g ~ bill_depth_mm, data = penguins)
2  typeof(example_mod)
```

```
[1] "list"
```

```r
1  length(example_mod$residuals);length(example_mod$coefficients)
```

```
[1] 342
```

```
[1] 2
```

# A Quick Aside on Naming Stuff

- Things we can never name stuff

```
 1  if
 2  else
 3  while
 4  function
 5  for
 6  TRUE
 7  FALSE
 8  NULL
 9  Inf
10  NaN
11  NA
```

There are more see this website for a more complete list

# A Quick Aside on Naming Stuff(cont)

Semi-reserved words

For simple things like assigning `c <- 4` and then doing `d <- c(1,2,3,4)` R will be able to distinguish between assign c the value of 4 and the `c` that calls `concatenate` which is way more important in R.

However it is generally a good idea, *unless you know what you are doing,* to avoid naming things that are functions in `R` because `R` will get confused.

```r
1  my_cool_fun <- function(x){
2    x <- x*5
3  return(x)
4  }
5
6  datas <- c(1:10)
7
8  my_cool_fun(datas)
```

```
[1]   5 10 15 20 25 30 35 40 45 50
```

```r
1  my_cool_fun[1]
```

```
Error in my_cool_fun[1]: object of type 'closure' is not subsettable
```

# How and What to Name Objects

The best practice is to use concise descriptive names

When loading in data typically I do `raw_my_dataset_name` and after data all of my cleaning I do `clean_my_dataset_name`

- Objects must start with a letter. But can contain letters, numbers, `_`, or `.`
  - snake_case_like_this_is_what_I_use
  - somePeopleUseCamelCase
  - some_People.are_Do_not.like_Convention

Example and Discussion provided in R for Data Sciency by Hadley Wickham

# Your Turn

- Create a vector from 1:100

- Create a character vector named hp with only the value of harry potter

- Find the length of each vector

- create a vector named pak to install "marginaleffects", "modelsummary"

04:00

# Navigating Objects in R

# Our Data

| species | island | bill_length_mm | bill_depth_mm | flipper_length_mm | b |
|---|---|---|---|---|---|
| Adelie | Torgersen | 39.1 | 18.7 | 181 | |
| Adelie | Torgersen | 39.5 | 17.4 | 186 | |
| Adelie | Torgersen | 40.3 | 18.0 | 195 | |
| Adelie | Torgersen | NA | NA | NA | |
| Adelie | Torgersen | 36.7 | 19.3 | 193 | |
| Adelie | Torgersen | 39.3 | 20.6 | 190 | |

# Indexing [ ]

- We can use column position to index objects.

- There are two slots we can use *rows* and *columns* in the brackets if we are using a dataframe like this.

- `object_name[row number, column number]`

- We can also subset our data by column position using `:` or `c(column 1, column 2)`

```
1  penguins[1,1]
```

| species |
| --- |
| Adelie |

```
1  penguins[1,1:2]
2  penguins[1,c(1,4)]
```

| species | island |
| --- | --- |
| Adelie | Torgersen |

| species | bill_depth_mm |
| --- | --- |
| Adelie | 18.7 |

# Indexing [ ] (cont)

- We can tell R what element of a list using a combo of `[ ]` and `[ [ ] ]`

```
1  my_list <- list(a = 1:4, b = "Hello World", c = data.frame(x = 1:3, y = 4:6))
```

```
1  my_list[[1]][2] ## get the first item in the list and the second element of that item
```

```
[1] 2
```
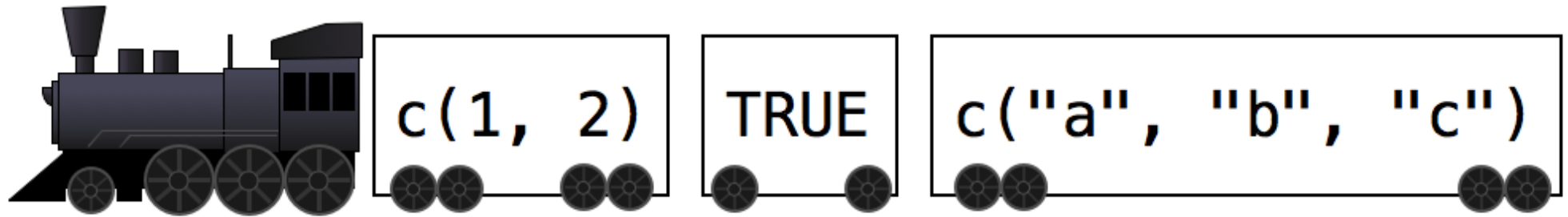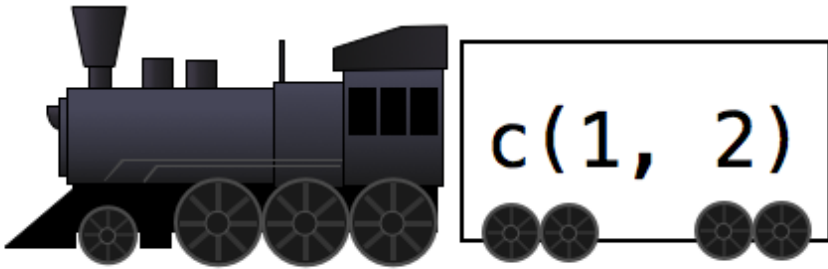
```
1  my_list[2]
```

```
$b
[1] "Hello World"
```

```
1  my_list[[3]][[1]]
```

```
[1] 1 2 3
```

# [ ] vs [ [ ] ]



c(1, 2)   TRUE   c("a", "b", "c")

lst

c(1, 2)

lst[1]          lst[[1]]

c(1, 2)

# Negative Indexing

- We can also exclude various elements using – and/or tests that I showed you earlier

```
1 penguins[,-1]
```

| island | bill_length_mm | bill_depth_mm | flipper_length_mm | body_mass_ |
|--------|---------------:|--------------:|------------------:|-----------:|
| Torgersen | 39.1 | 18.7 | 181 | 375 |
| Torgersen | 39.5 | 17.4 | 186 | 380 |
| Torgersen | 40.3 | 18.0 | 195 | 325 |
| Torgersen | NA | NA | NA | N |
| Torgersen | 36.7 | 19.3 | 193 | 345 |
| Torgersen | 39.3 | 20.6 | 190 | 365 |

# Negative Indexing(cont)

- We can use – or : as well to subset stuff

```
1  penguins[,-(1:4)]
```

| flipper_length_mm | body_mass_g | se |
|---:|---:|---|
| 181 | 3750 | ma |
| 186 | 3800 | fe |
| 195 | 3250 | fe |
| NA | NA | NA |
| 193 | 3450 | fe |
| 190 | 3650 | ma |

```
1  penguins[,-c(2,3,5,8)]
```

| species | bill_depth_mm | body_mass_ |
|---|---:|---:|
| Adelie | 18.7 | 375 |
| Adelie | 17.4 | 380 |
| Adelie | 18.0 | 325 |
| Adelie | NA | N |
| Adelie | 19.3 | 345 |
| Adelie | 20.6 | 365 |

# Subsetting By Tests

```
1  penguins[penguins["sex"] == "female", c("species", "sex")]
```

| species | sex |
|---------|--------|
| Adelie | female |
| Adelie | female |
| NA | NA |
| Adelie | female |
| Adelie | female |
| NA | NA |
| NA | NA |
| NA | NA |
| NA | NA |
| Adelie | female |

# $ Indexing

A really useful way of indexing in ℝ is referencing stuff by name rather than position. - The way we do this is throught the $

```
1 my_list$a
```

```
[1] 1 2 3 4
```

```
1 my_list$b
```

```
[1] "Hello World"
```

```
1 my_list$c
```

```
  x y
1 1 4
2 2 5
3 3 6
```

# Indexing(cont)

```r
1  my_list[[3]][[2]] ## these are just returning the same thing
```

```
[1] 4 5 6
```

```r
1  my_list$c$y
```

```
[1] 4 5 6
```

# $ in action

This will just subset things

```
1  penguins[penguins$species == "Gentoo", c("species", "island", "bill_length_mm")]
```

| species | island | bill_length_mm |
|---------|--------|----------------|
| Gentoo  | Biscoe | 46.1 |
| Gentoo  | Biscoe | 50.0 |
| Gentoo  | Biscoe | 48.7 |
| Gentoo  | Biscoe | 50.0 |
| Gentoo  | Biscoe | 47.6 |
| Gentoo  | Biscoe | 46.5 |
| Gentoo  | Biscoe | 45.4 |
| Gentoo  | Biscoe | 46.7 |
| Gentoo  | Biscoe | 43.3 |
| Gentoo  | Biscoe | 46.8 |

# Comparing what we know how to do

## Tidyverse

```
1 penguins |>
2 select(species, island, sex)
```

| species | island | sex |
|---------|-----------|--------|
| Adelie | Torgersen | male |
| Adelie | Torgersen | female |
| Adelie | Torgersen | female |
| Adelie | Torgersen | NA |
| Adelie | Torgersen | female |

## Base R

```
1 penguins[, c("species", "island", "sex")]
```

| species | island | sex |
|---------|-----------|--------|
| Adelie | Torgersen | male |
| Adelie | Torgersen | female |
| Adelie | Torgersen | female |
| Adelie | Torgersen | NA |
| Adelie | Torgersen | female |

# Sometimes it is just quicker

```r
1  penguins_base$range_body_mass <- max(penguins_base
2      -
3  penguins_base$bill_ratio  <- penguins_base$bill_le
4
5  mean(penguins_base$body_mass_g, na.rm = TRUE)
```

```
[1] 4201.754
```

```r
1  penguins <- penguins |>
2  mutate(range_body_mass = max(body_mass_g, na.rm =
3          bill_ration = bill_length_mm/bill_depth_mm)
4
5  summarise(penguins, mean(body_mass_g, na.rm = TRUE
```

```
# A tibble: 1 × 1
  `mean(body_mass_g, na.rm =
TRUE)`

<dbl>
1
4202.
```

# Sometimes the Original is Just as Good as the Wrapper

```
1  data("starwars")
2
3  filter(starwars, str_detect(eye_color, "blu"))
```

| name | eye_color |
| --- | --- |
| Luke Skywalker | blue |
| Owen Lars | blue |
| Beru Whitesun lars | blue |
| Obi-Wan Kenobi | blue-gray |
| Anakin Skywalker | blue |
| Wilhuff Tarkin | blue |
| Chewbacca | blue |
| Jek Tono Porkins | blue |

```
1  starwars[grepl("blu",starwars$eye_color),]
```

| name | eye_color |
| --- | --- |
| Luke Skywalker | blue |
| Owen Lars | blue |
| Beru Whitesun lars | blue |
| Obi-Wan Kenobi | blue-gray |
| Anakin Skywalker | blue |
| Wilhuff Tarkin | blue |
| Chewbacca | blue |
| Jek Tono Porkins | blue |

# Finding Help

- Asking for help in R is easy the most common ways are
  `help(thingineedhelpwith)` and `?thingineedhelpwith`

```
1  ?grepl
```

- `?thingineedhelpwith` is probably the most common because it requires less typing.

- Base and Tidy functions differ in many ways other than naming conventions

# Finding Help

## Extract highest and lowest grades from a data frame

**Description**

grading is just a package that returns grades for exams from highest to lowest

**Usage**

```
grading( ... , data = NULL)
```

**Arguments**

...
    pass of starts_with() since that will make your life easiers

data
    is a data.frame that you need graded

**Details**

grading requires some exam data

---

*Cutting corners to meet arbitrary management deadlines*

*Essential*

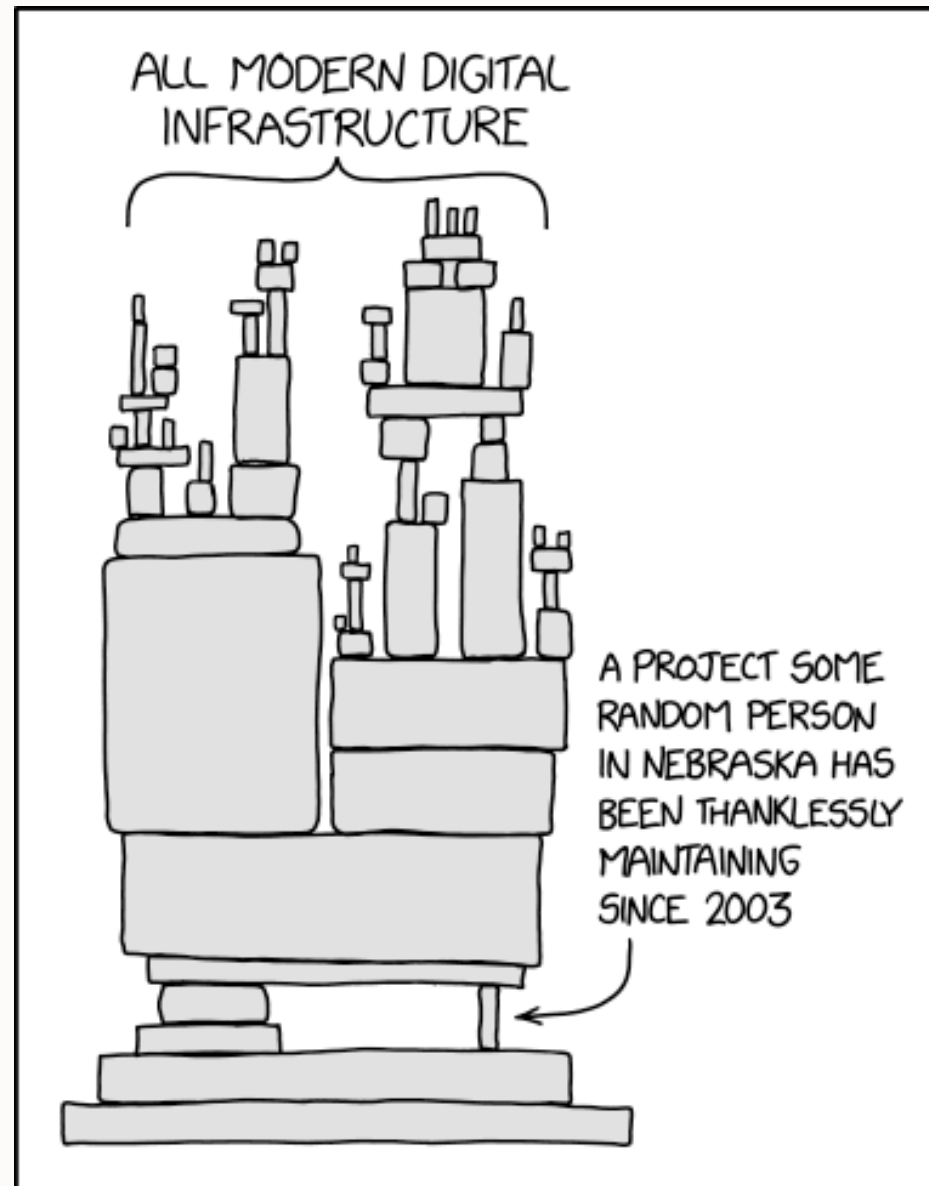## Copying and Pasting from Stack Overflow

O'REILLY®

*The Practical Developer*
*@ThePracticalDev*

# Your Turn

- Find the minimum value of `bill_length_mm`

- Find the maximum value of `body_mass_g`

- Subset the penguins data any way you want using column position or `$`

- Assign each of them to an object

- Create a vector from 1:10 index that vector using `[ ]` to return 2 and 4

# The Tidyverse issue

# Learning to Live With Each other

```r
penguins$big_peng <-  dplyr::case_when(penguins$body_mass_g > mean(penguins$body_mass_g, na.rm = TRUE) ~ "Big
     penguins$body_mass_g < mean(penguins$body_mass_g, na.rm = TRUE) ~ "Smol Penguin",
     penguins$body_mass_g == mean(penguins$body_mass_g, na.rm = TRUE) ~ "Average Penguin")

penguins$body_mass_g[is.na(penguins$body_mass_g)] <- 0

penguins |>
  select(body_mass_g, big_peng)
```

```
# A tibble: 344 × 2
  body_mass_g big_peng
        <dbl> <chr>
1        3750 Smol Penguin
2        3800 Smol Penguin
3        3250 Smol Penguin
4           0 <NA>
5        3450 Smol Penguin
6        3650 Smol Penguin
7        3625 Smol Penguin
8        4675 Big Penguin
9        3475 Smol Penguin
```

# Learning to Live with Each Other

- Lots of stuff is repetitive

- Repetition isn't necessarily bad but it can easily lead to mistakes

```
1  penguins |>
2  drop_na() |>
3  mutate(body_mass_g = body_mass_g - min(body_mass_g, na.rm = TRUE) /
4      (max(body_mass_g, na.rm = TRUE) - min(body_mass_g, na.rm = TRUE)),
5      flipper_length_mm = flipper_length_mm - min(flipper_length_mm, na.rm = TRUE) /
6      (max(flipper_length_mm, na.rm = TRUE) - min(flipper_length_mm, na.rm = TRUE)),
7      bill_length_mm = bill_length_mm - min(bill_length_mm, na.rm = TRUE) /
8      (max(bill_length_mm, na.rm = TRUE) - min(flipper_length_mm, na.rm = TRUE)))
```

# One way this helps us

```r
# we can rewrite this code pretty easily and iterate over the entire dataset
rescale <- function(x){
    rng <- range(x, na.rm = TRUE, finite = TRUE)

  x-rng[1]/(rng[2] - rng[1])

}

penguins |>
mutate(across(where(is.numeric), \(x) rescale(x)))
```

| species | island | bill_length_mm | bill_depth_mm | flipper_length_mm | b |
|---|---|---|---|---|---|
| Adelie | Torgersen | 37.93273 | 17.14048 | 178.0847 | |
| Adelie | Torgersen | 38.33273 | 15.84048 | 183.0847 | |
| Adelie | Torgersen | 39.13273 | 16.44048 | 192.0847 | |

| species | island | bill_length_mm | bill_depth_mm | flipper_length_mm | b |
|---------|--------|---------------:|--------------:|------------------:|---|
| Adelie | Torgersen | NA | NA | NA | |
| Adelie | Torgersen | 35.53273 | 17.74048 | 190.0847 | |

# Reading in Data Gets easier

```r
1   rm(list = ls())
2   penguins <- palmerpenguins::penguins
3   starwars <- dplyr::starwars
4   data("mpg")
5   data("mtcars")
6
7   data_names = c("mpg", "penguins", "starwars", "mtcars")
8
9   for(i in 1:length(data_names)) {
10    readr::write_csv(get(data_names[i]),
11            paste0("data/",
12                    data_names[i],
13                    ".csv"))
14  }
15  my_files <- list.files(path = "data/",pattern = "*.csv", full.names = TRUE)
16
17  # Further arguments to read.csv can be passed in ...
```

# Which Gives Us

```
[1] "all_csv"    "data_names" "i"          "mpg"        "mtcars"
[6] "my_files"   "penguins"   "starwars"
```