# Day 2, Part 1: R Basics

Brennan Terhune-Cotter and Matt Dye

https://github.com/brennangitsit/2023\_IAM3\_R

#### Review

- Using R for reproducible analyses :)
- How to use RStudio and its features
- How to import data

## Agenda

- 1. The R Environment
- 2. Objects
- 3. Operators
- 4. How to Write Good Code

# The R Environment

#### The R Environment

- *R* is a **free** software environment for statistical computing & plotting
  - R is open-source; anyone can contribute to R.
  - Thousands of free *packages* have been developed by different contributors.
- Like most other open-source languages, you can work with R via many user *interfaces*.
  - However, it's best to use an IDE (integrated development environment).
  - RStudio is the most common IDE for R. It gives us tools to more effectively use R.

#### Base R & CRAN

- *Base R* is the basic software containing the R programming language.
  - Frequently updated; most recent is *Already Tomorrow* (04-21-2023)
- The Comprehensive R Archive Network (CRAN) has:
  - Mirrors for downloading R code and documentation
  - A repository of *packages* (now has ~19,000 packages!)

#### Packages

- Themed collection of functions and datasets with a cohesive purpose
- Extend basic R functionality; essential to using R

## Installing and Using Packages

Install packages once per machine; use the following commands in your console (or click prompts to "install" when they pop up in RStudio)

(package names are in quotes)

Load packages in every session; use the following commands in your scripts

(package names are NOT in quotes)

1 library(tidyverse)

<sup>1</sup> install.packages("tidyverse")

## Installing and Using Packages

• When a package is loaded, you can use any function in that package just by calling the function name

```
1 library(readxl)
2 read_xlsx(...)
```

• However, you can also call any function from an installed package without loading it by adding the package name:

```
1 readxl::read_xlsx(...)
```

- This is useful if:
  - You just want to use one function from a package
  - The function is *masked* by another function from another package with the same name

# **R** Files

#### **R** Projects

- An R project (.Rproj) is a file which saves your workspace
- It is useful to have an .Rproj file for each project you have
- When you open the .Rproj, everything will look like the last time you saved it.
- Instead of opening RStudio or individual R files (like any other application):
  - 1. Open your .Rproj
  - 2. Open files from that . Rproj  $\P$



#### **R** Scripts

#### library(readxl) library(dplyr) getwd() setwd("data/") cpi <- read\_xlsx("cpi\_raw.xlsx") %>% janitor::clean\_names() 12 cpi\_clean <- cpi %>% filter(!is.na(consumer\_price\_index\_item)) %>% filter(!is.na(annual\_2021)) %>% $rename(predicted_2023 = x11,$ hist\_avg\_2003\_2022 = x20\_year\_historical\_average\_2003\_2022, item = consumer\_price\_index\_item) %>% select(item, annual\_2020, annual\_2021, annual\_2022, predicted\_2023, hist\_avg\_2 saveRDS(cpi\_clean, "cpi.rds") tidy\_cpi <- cpi %>% rename(annual\_2023 = predicted\_2023) %>% pivot\_longer(annual\_2020:annual\_2023, names\_to = "year", names\_prefix = "annual\_", values\_to = "increase") %>% select(item,year,increase,everything()) %>% saveRDS("cpi\_tidy.rds") eruptions <- read\_xlsx("holocene\_eruptions.xlsx") %>% janitor::clean\_names() 35 • extract\_year <- function(df, col\_name) {</pre> library(stringr) df %>% mutate(year = str\_extract(last\_known\_eruption, "\\d+")) %>% mutate(year = as.numeric(year)) %>% mutate(year = if\_else(str\_detect(last\_known\_eruption, "BCE"), year \* -1, yea select(volcano\_number,volcano\_name,country,last\_known\_eruption,year,everythi

- *Scripts* have *code* which can be executed
- R scripts end with an '.R' extension
- Everything in an R script will be executed unless it is *commented* out with #

#### **R** Markdown



- RStudio also provides R Markdown/Notebook documents (".Rmd")
- Rmd files are *text files* that contain executable *code chunks*
- More readable & great for sharing your data analysis because...

#### **R** Markdown

#### Importing Data

library(tidyverse)

Next, we use the read\_csv() function to pull in data from our survey. We tell R to put that data in a data frame called climate\_data. If it finds any blank cells, then those are NA data. NA values are treated specially in R; we will cover how to work with NA values later.

Code 🕶

Hide

Hide

After running this command, a new 'object' called climate\_data appears in the Environment pane in the top right corner of R Studio. Click on it and the dataframe opens in a new window and can be inspected. This is equivalent to typing View(climate\_data) in the Console below.

Always be careful with dates. Looking at our CSV file, we can see that it is in M/DD/YY HH:MM format. In R terminology, this is %m%d%y %H%M. It is best to read dates into a data frame in ISO8601 format, which is %Y-%m-%d %H:%M:%S. In tidyverse this is refereed to as "datetime" format which has the data type <dttm>. We use the col\_types argument in read\_csv to do this.

climate_data <- read_csv(*~/Documents/GitHub/2023_IAM3_R/data/climate-survey-data.csv*,	
na = "",	
<pre>col_types = cols(StartDate = col_datetime("%m/%d/%y %H:%M"),</pre>	
EndDate = col_datetime("%m/%d/%y %H:%M"),	
RecordedDate = col_datetime("%m/%d/%y %H:%M"))	
)	

Our next step is to rename the variables so that they conform to good R practice. Note that by convention we put (what we will be creating / what we want) to the left of the expression.



- .Rmd files are *rendered* into readable .html files
- Those files show your commentary, code, and output (including plots!)
- Two types of .Rmd files:
  - *R Markdown* files execute all code chunks upon "knitting" (rendering)
  - R Notebook files only include output of code you've already executed
    - Newer format & more interactive
    - Act as a lab "notebook"



#### **Objects**

- Basic building blocks of your code and data!
- Created using assignment operators <- or =
- Can be manipulated using their names

```
1 object <- 2
2 object
[1] 2
[1] anotherobject <- c(4,5,6)
2 anotherobject
[1] 4 5 6
1 object + anotherobject</pre>
```

[1] 6 7 8

#### **Objects**

○1. Open 13\_basics.Rmd

2. Read Assigning Variables

## **Data Types**

Each object has a **type**. R has several data types. We only need to know 4:

- Numeric (or double)
  - Integer is another type which is just numeric without decimal points
- Character (or string)
  - Characters or sets of characters (including numbers)
  - Strings must be surrounded by single (''') or double ('''') quotes
- Boolean (or logical)
  - Boolean variables can only take two values: TRUE/T or FALSE/F

```
1 numeric <- 2
2 character <- 'a'
3 another_character <- '2'
4 string <- "this is stringy"
5 boolean <- TRUE
6 another logical <- F</pre>
```



**Q** Read **Data Types** 

#### **Factors**

- The 4th data type to know is **factors**
- Factors are used to represent categorical data
  - Data that can take on a limited number of distinct values (these values are referred to as levels).
  - If you have a factor representing colors, its levels might be "red", "blue", and "green".
- Factors can be **unordered** (the default), where no level is considered "greater" than any other, or they can be **ordered**, where the levels have a specific sequence.
  - A factor representing the t-shirt sizes "small", "medium", and "large" could be an ordered factor.

## Data Types: NA

<ul> <li>R treats missing values with a special type called NA</li> </ul>	no_NA	has_NA
<ul> <li>This is NOT like typing "NA" as a string!</li> </ul>	а	d
<ul> <li>NA values require special functions</li> </ul>	NA	NA
<ul> <li>They can be really annoying to deal with but helps you</li> </ul>	b	NA
handle missing data appropriately!	NA	е

#### **Vectors**

- Essentially one "column" of data (one-dimensional array)
- In R, a single number is technically a vector of length 1
- All values in a vector must be the same data type!

```
1 num_vector <- c(1,2,3)
```

```
2 char_vector <- c("hi","what's","up")</pre>
```

- 3 boolean\_vector <- c(T,F,TRUE)</pre>
- 4 length(num\_vector)

#### [1] 3

```
1 number <- 4
```

2 length(number)

#### [1] 1

#### **Matrices**

- A matrix is a two-dimensional array
- All values in a matrix must be the same data type!

```
1 data <- c(1,2,3,4,5,6,7,8,9)
 2 matrix <- matrix (data, nrow=3, ncol=3)</pre>
  3 matrix
     [,1] [,2] [,3]
[1,]
        1
              4
                   7
[2,]
              5
                   8
        2
[3,]
        3
              6
                   9
  1 matrix[2,3]
[1] 8
```

## Data Frames (dfs)

#### 🖓 Tip

This is what we will be using 99% of the time when working with data in R!

- A two-dimensional data structure (rows and columns)
- Essentially a collection of named vectors (columns):
  - All data within a column must be the same type
  - BUT different columns within a df can have different data types!
  - And columns in a df have names

```
1 dataset <- data.frame(num_vector, char_vector, boolean_vector)
2 dataset</pre>
```

	num_vector	char_vector	boolean_vector
1	1	hi	TRUE
2	2	what's	FALSE
3	3	up	TRUE

#### **Data Frames**

our data in Excel worksheets == dataframes (with hardwired variable names)

	А	В	C
1	ID	Age	Score
2	01	22	5
3	02	56	6
4	03	103	8
5	04	4	11
6	05	67	12
7	06	31	2
8	07	-40	
a			

	ID <sup>‡</sup>	Age 🍦	Score <sup>‡</sup>
1	01	22	5
2	02	56	6
3	03	103	8
4	04	4	11
5	05	67	12
6	06	31	2
7	07	-40	NA

#### **Functions**

- Functions are blocks of code which perform a specific task
- Functions take *arguments* as input
- You will be using ready-made functions to manipulate your data (from base R or from packages)
- Near the end of the workshop, you can learn to create your own functions
  - Important for simpler, more organized & readable code

1 functionname(argument1, argument2, ...)



#### **Common Operators**

- Assignment operators: <- (option -)
- Relational operators: ==, !=, >, <, >=, <=
- Creating a series of numbers: :

```
1 y <- 1:10
2 print(y) # bonus: what kind of object is y?
[1] 1 2 3 4 5 6 7 8 9 10</pre>
```

#### Read Operators and NA

## How to Write Good Code

#### Commenting

- Commenting is your *inline documentation* of your code and analysis
- Especially as a beginning coder, there is no such thing as too little commenting
- Comments should:
  - 1. explain what your code is doing
  - 2. explain decisions you made and why
  - 3. not repeat the code, but clarify & contextualize it

## Naming Variables and Objects

How you name variables and objects can make life much easier for you.

- Use long & descriptive variable or object names if you have to. Text is cheap, brain capacity is not.
  - Which dataframe name is clearer?

```
1 df3
2 average_EEG_response_times
```

- Variables and objects should never have spaces; use underscores instead
  - Names with spaces must be surrounded by ``` every time you call them, which is super annoying
- Variable names should be lowercase

## **Coding Style**

- Don't use run-on code lines; most functions should start on a new line.
- Use blank lines often to separate code blocks! You can't have too many blank lines.
- Add spaces around operators: + == < != <- etc.
- Add spaces after comments like in English.

#### Pseudocode! (your salvation!)

- What do you do if you're not sure how to write something in code?
- You write *pseudocode* first, and gradually change it into code.

# Pseudocode 1 For data frame 'people', 2 Find rows where 'age' is greater than 20, 3 From these rows, calculate the mean of 'age' 1 With data frame 'people', 2 Use filter() to select rows where 'age' > 20, 3 Use summarise() with mean() to 4 calculate mean 'age' of these rows.

#### Code

- library(dplyr)
- 2
- 3 people %>%
- 4 filter(age > 20) %>%
- 5 summarise(mean\_age = mean(age))

## Debugging

- Overwriting objects (or not!)
  - When you assign the same name to an object, you overwrite that object.
  - Assigning different names allows you to look at the object in between your steps so you can figure out where you went wrong.
- Rerunning code
  - You should make a habit of rerunning your entire code with a clean dataset pretty often.
  - This helps you catch mistakes early on.

## Other tips for coding in R

- Delete old objects you no longer need with rm(). This helps keep your environment clean.
- If you need to quote something, highlight it and press " or '. This also works with (.
- Use **sectioning comments** (# Section title -----) which allow you to "minimize" sections.