

A Walk on the Side

an introduction to R for data analysis



GW Libraries Workshop
Spring 2024

go.gwu.edu/rworkshop

FAQ

Q: Will you sign my form for Professional Enhancement hours?

A: Yes, email me!

Q: Can I get a copy of your R code?

A: Yes, email me

Q: Will this workshop be recorded?

A: No, so hang on for the ride!



Logistics

- Schedule

9:30- 2:15 with a ~1 hr break for lunch



Upcoming R workshops

- Feb. 1 (Thurs., 9:30am-2:15pm)

This workshop, AGAIN! ← *tell your friends!*

- Feb. 6 (Tues., 1-3:30pm)

Farther into R: *More R for Data Analysis*





Goals



Learning Objectives

[Hopefully] You will learn how to do some of the following:

- Set up your laptop with R & RStudio (done!)
- Write and run an R program in RStudio
- Use variables of different types in R
- Use vectors and data frames in R to represent data
- Import & export data files
- "Wrangle" data in R
- Explore data in R with basic statistics and data visualizations
- Learn how to look for help to overcome obstacles



Agenda

- About R and RStudio
- Along the way: How to get help
- Hands-on:
 - variables
 - logical expressions
 - values, vectors, and data frames
 - R Studio projects
 - reading in data
 - exploring data
 - data wrangling:
cleaning and reshaping
 - data visualization
 - data analysis
 - functions
 - R Markdown / reports
- Resources for further learning



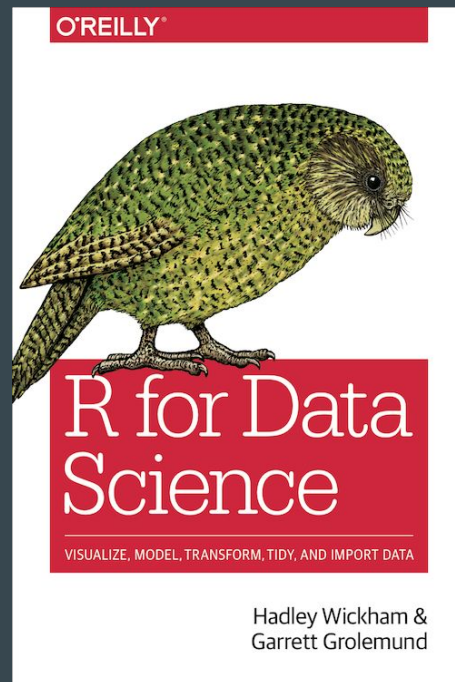
Acknowledgments



Teaching basic lab skills
for research computing



r-tutor.com



r4ds.had.co.nz



Workshop Housekeeping



Ask questions! Either via voice or chat

Use chat to help each other out

If something is confusing in the workshop, let us know.

About R

- Free/Open source
- Cross-platform (Mac, Windows, Linux)
- For statistical computing (and data visualization)
- CRAN - r-project.org
 - [R packages](#)
 - [R journal](#)

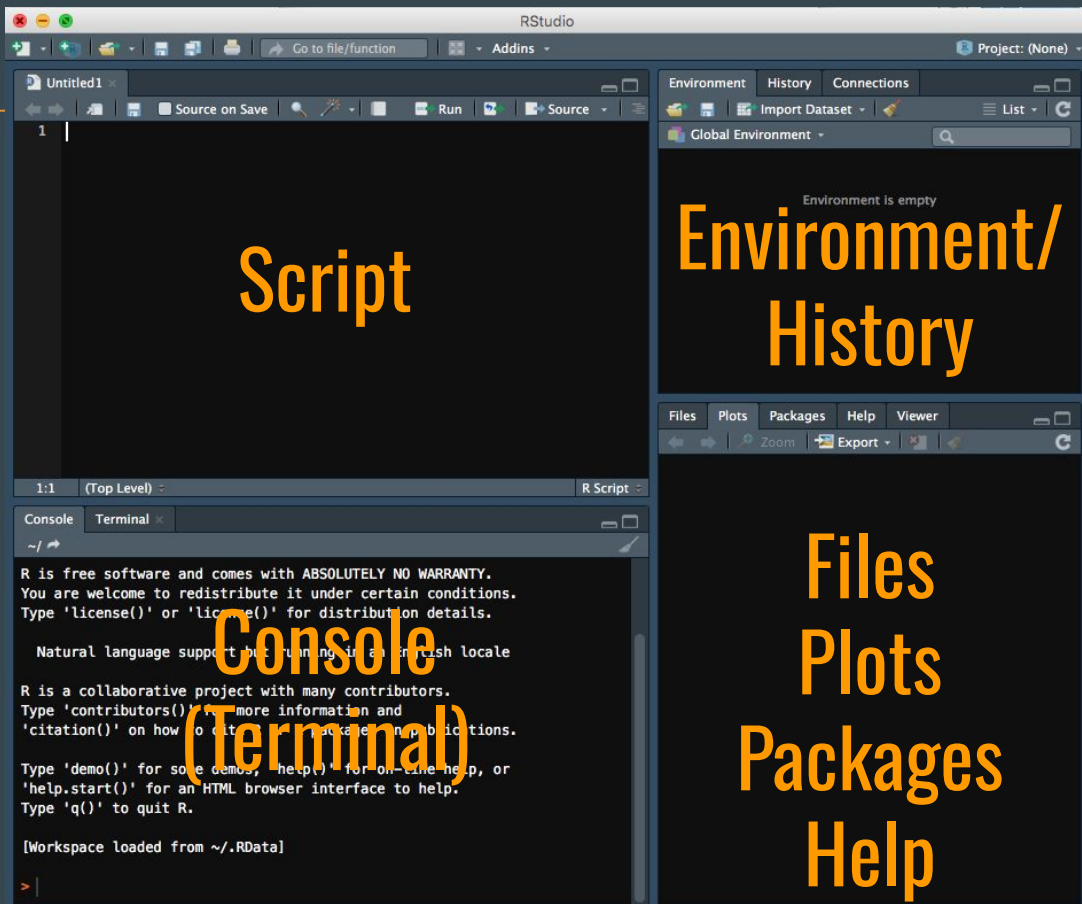




Reasons researchers prefer R

- Scripted language (vs. point/click)
- Features built around working with data
- Reproducibility
- Interdisciplinary
- Extensible
- Beautiful data visualization
- RStudio (Posit) is a well-liked R development app
- Community - RStudio Community, Stack Overflow

R Studio



A WALK ON THE R SIDE



Variables/Objects

"Binding" data to a named object/variable allows you to store data in memory and access it later.

```
x <- 5
```

```
y <- c("Washington", "Chicago", "Washington", "Boston")
```

```
z <- data.frame(pt_id = c("A001", "B204"), bpm = c(60, 72))
```





Variables

- Try using R as a "calculator" in the Console
 - Try some mathematical functions, too
- Create some variables
 - variable naming
 - `<-` for assigning values to variables (Option - on Mac, Alt - on Win)
 - numeric, character, logical
 - Watch the Environment pane!
 - `typeof()`
 - Coercion w/ `as.integer`, `as.character`, `as.logical`, `as...`

Logical Expressions

- Operators include:
==, <, >, ! (not), & (and), | (or), etc.





Basic Data Structures

Atomic Vector

10.2

Vector

1	10.2
2	11.3
3	11.5
4	12.0

Data Frame

	time	temp	boiling
1	51	10.2	FALSE
2	58	11.3	FALSE
3	63	11.5	FALSE
4	70	12.0	TRUE



Vectors

Vectors

- A vector is
 - A sequence of data elements (components) all of the same type.
- Create vectors with `c()` (short for "combine")



The image shows a screenshot of the RStudio integrated development environment. The main window displays an R script with the following code:

```
1 library('tidyverse')
2 gapminder <- read_csv('data/gapminder.csv')
3
4 by_year <- gapminder %>%
5   group_by(year) %>%
6   summarize(weighted_avg_lifeExp = sum(pop*lifeExp)/sum(pop))
7
8 # Plot the data (scatterplot)
9 plot(y = by_year$weighted_avg_lifeExp, x = by_year$year, col='blue')
10 # Build a linear regression model
11 mod = lm(data = by_year, weighted_avg_lifeExp ~ year)
12 # Plot the line
13 abline(mod)
14
15 # or using ggplot2:
16 ggplot(data = gapminder, aes(x=year, y=lifeExp, group=continent, color=continent)) +
17   geom_point() +
18   # ...
19
20 # ...
```

The console window shows the output of the code, including the structure of the 'gapminder' data frame and the results of the 'summarize' function applied to the 'by_year' data frame.

On the right side, the 'Environment' pane shows the objects in the global environment, including 'df', 'gapminder', 'housedata', 'lemod', 'mod', 'mx', and 'mx2'. The 'Values' pane shows the values of the 'primes' and 'testnum' objects.

A large, semi-transparent text overlay is centered on the image, reading: "Let's pause to explore some useful tabs in RStudio".



Data Frames



Data Frames

- A **data.frame** stores a data table
- Comprised of **vectors** of equal length.
Vectors become **columns**.
- Columns and rows can have names.
- **tibble** (from the tibble package) has some advantages over **data.frame**



A brief word on **list** and **matrix**



Projects in RStudio

Projects in RStudio

Recommendations:

- Use [Github for] **version control!**
- Create **folders** to keep things organized





It's time to **import** some data!



Data Importing

- Prepare data as "tidy"
 - rectangular
 - one table per file
 - rows are observations, columns are variables
- Formats: CSV, TSV, Excel, Fixed-Width, JSON... and with the right packages: Stata, SPSS, SAS... (using **rio** or **haven**)
- A word about "big data" (consider **data.table**)



R Packages

Installing and loading R packages

- `install.packages('mypackage')`
- `library(mypackage)`

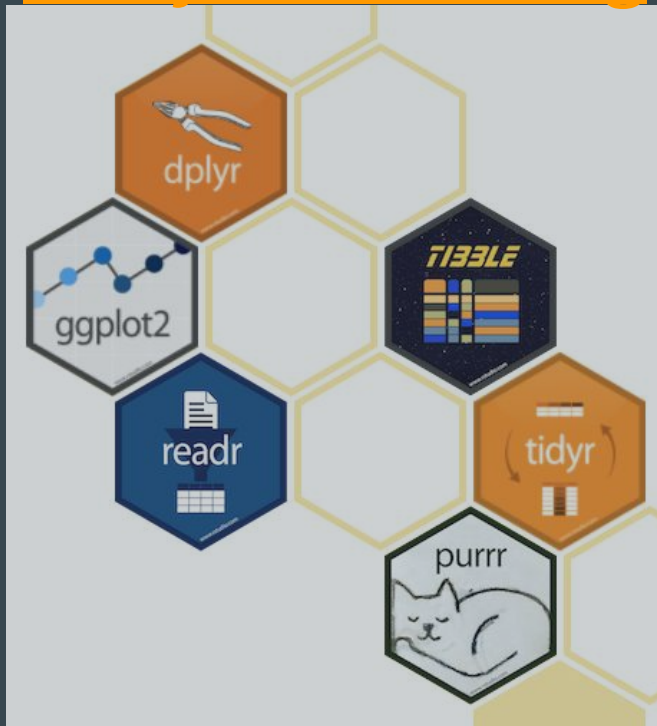




Tidyverse Core Packages

tidyverse.org

- ggplot2 - graphics
- dplyr - data manipulation
- tidyr - tidying data
- readr - reading in data
- tibble - modern data frame
- purrr - functional programming



Other often-used R packages

Loading in various data file types ♦ haven, readxl

Mapping ♦ rgdal, tmap, leaflet

Analyzing 2D and 3D shapes ♦ geomorph

Genomic data ♦ bioconductor

Cluster analyses ♦ cluster

Time series data ♦ forecast

Text mining ♦ qdap, sentimentr, tidytext

graph/network analysis ♦ igraph, sna

Interactive web visualizations ♦ shiny

Web scraping ♦ rvest



Exploring Data

- head, tail
- subsetting
- slicing and dicing





Data Wrangling

[flickr.com/photos/thewomensmuseum/3697075917](https://www.flickr.com/photos/thewomensmuseum/3697075917)

Data Transformation using the dplyr package

- `select()` # keep only certain columns
- `filter()` # keep only certain rows
- `mutate()` # add/modify variables
- `group_by() %>% summarize()`
compute summary statistics per group
- `arrange()` # order by a variable
- `dropna()` # drop rows with NAs in specified vars.

You will want to use a "pipe": `%>%`
(shortcut: **control-shift-M**)





Joining with dplyr

"Merge" tables together

- `left_join()`
- `right_join()`
- ...

Data Tidying/Reshaping with tidyr

- `pivot_wider()`
- `pivot_longer()`
- ...



Data Visualization with "base R" and `ggplot`



Data Analysis



Functions



R Markdown



R Markdown

- A format for writing reproducible, dynamic reports with R (as HTML, PDF, MS Word, and more)
- rmarkdown.rstudio.com
- # Header 1
Header 2
Italic ****bold****
- Insert R code directly into your document

```
```{r setup}
your R code goes here
```
```
- Include LaTeX code with \$ or \$\$



R Shiny



Parting thoughts



Recommended practices

- Use Projects in RStudio
 - Set up folders
- Use tidyverse packages (dplyr, tidyr, etc.) to wrangle your data
- Leave raw data raw
- 🪲 Empty out your variables, then make sure your script runs from the top
- Learn by finding and using working examples



Some Handy R Links

NEW for 2024!! R "libguide"



Only the **best** R links:

libguides.gwu.edu/Rstats

Thanks!

Dan Kerchner kerchner@gwu.edu

These slides: go.gwu.edu/rworkshop

Statistics focused (+ R/Python/SAS/etc.) appointments
w/graduate student consultants: go.gwu.edu/dataconsulting

Appointments with me: calendly.com/kerchner

Coding consultations (**R**, Python, HTML/CSS, etc.):
calendly.com/gwul-coding

