# R for Data

**How to Get Started**

## How do you currently use data?
## What is your interest in R?

# Why R is Good

Why would I want to use it?

Is it worth my time?

- Free
  - Can be used anywhere, for any purpose

- Syntax = Reproducibility
  - Statistical software also has syntax

- Data Notebook Capable
  - Easily write data documents

- Unlimited Capabilities
  - Easy to do complex things

GEORGE MASON UNIVERSITY

# Why R is Hard

- Text-based computer language
  - must know what you want
  - must use exact syntax

- Constantly changing, huge ecosystem
  - albeit improving

- Many ways to do the same thing
  - Many people add functionality
  - Too may tutorials

GEORGE MASON UNIVERSITY

# Changes over time

- **RStudio → Posit**

- **R Markdown → Quarto Document**

- **%>% → |>**

- **plyr → dplyr**

- **reshape2 → tidyr**

# Low on Time?

Utilize Graphical User Interfaces

Start learning basic programming concepts

Use it just for specific tasks

# Crash Course

If you know these already, then you are in great shape.

Programmers often find R confusing / harder.

Easier to teach R to non-programmers

GEORGE MASON UNIVERSITY

# Principles: Console vs Script

**>** prompt

**+** waiting for more

**[1]** returned element

Code separate, saves in file

"Run" or "Execute" code

# is a comment

```
Type 'de
'help.st
Type 'q(

> 3+2
[1] 5
> 3-
+ 2
[1] 1
> |
```

```
  1  setwd("C:
  2
  3  library(t
  4  test <- r
  5  train <-
  6
  7  library(s
  8  dfSummary
  9
 10  library(G
 11  train %>%
 12     select(
 13     ggpairs
 14
```

GEORGE
MASON
UNIVERSITY

# Principles

## Functions

Arguments

```
function( )
```

Value

```
sum(2,4)
countif(A2:A50,">5")
read.table(mydata, header=TRUE, sep=",")
```

## Objects

*Programmers call these variables. But columns in a data table are variables to statisticians.*

- Numbers
- Characters

*Data Types*

- Vectors
- Datasets

*Data Structures*

```
people <- 9
school <- "GMU"
```

*Assignment Operator*

GEORGE MASON UNIVERSITY

Data Types

# Data Structures

|  | Homogeneous (same Type) | Heterogeneous (different Types) |
|---|---|---|
| 1 Dimension | **Vector** | **List** |
| 2 Dimensions | **Matrix** | **Data Frame** |

# Principles: Data Tables & Tidy Data



variables            observations            values

https://r4ds.had.co.nz/tidy-data.html

# Programming Concepts that Help

| | | |
|---|---|---|
| **Data Types** (e.g., Numeric, String) | **Data Structures** (e.g., Lists, Dictionaries) | **Variables** and **Objects** |
| **Conditional** Statements (If Statements) | **Functions**, **Methods**, and **Modules** | **Working Directory** and **File Paths** |

GEORGE MASON UNIVERSITY

# Quick Start: Posit Primers & learnr Tutorials

## Posit Primers

### The Basics
Start here to learn the skills that you will rely on in every analysis (and every primer that follows): how to inspect, visualize, subset, and transform your data, as well as how to run code.

### Work with Data
Learn the most important data handling skills in R: how to extract values from a table, subset tables, calculate summary statistics, and derive new variables.

### Visualize Data
Learn how to use ggplot2 to make any type of plot with your data. Then learn the best ways to visualize patterns within values and relationships between variables.

### Tidy Your Data
Unlock the tidyverse by learning how to make and use tidy data, the data format designed for R.

### Iterate
Master a core programming paradigm with the purrr package: for each ___ do ___.

### Write Functions
Functions are the key to programming in R. This primer will teach you how to write and use your own reusable functions.

---

**Environment | History | Connections | Tutorial**

### Data basics
*learnr: ex-data-basics*

Learn about the base data types in R. Explore R's data frames, and learn how to interact with data frames and their columns.

**Start Tutorial ▶**

### Filter observations
*learnr: ex-data-filter*

Learn how to filter observations in a data frame. Use

**Start Tutorial ▶**

---

https://posit.cloud/learn/primers

https://rstudio.github.io/learnr/articles/examples.html

GEORGE MASON UNIVERSITY

# RStudio Education

https://education.rstudio.com/

**FOR LEARNERS**

## Expand your knowledge

Dive deeper into our popular packages like tidyverse and Shiny, with resources for beginner, intermediate, and expert-level R learners.

**LEARN MORE**

**FOR TEACHERS**

## Explore our resources

Use our materials and evidence-based teaching practices to teach data science using R and RStudio's products.

**LEARN MORE**

### Beginners

*Get started with the Tidyverse and R Markdown.* No one starting point will serve all beginners, but here are 6 ways to begin learning R. Read more ...

### Intermediates

*Expand your R skills.* Here are some common areas that people who already have some experience in R find particularly rewarding to learn. Read more ...

### Experts

*Go deep.* Learning some topics in depth will both help you develop better code and share it more effectively with others. Read more ...

# What do I need to get started?

**THE SOFTWARE**

**AN INTERFACE**

**PACKAGES & FUNCTIONS**

https://learnr-examples.shinyapps.io/ex-setup-r/

# The Software

- https://cran.r-project.org/

# An Interface – RStudio

**Desktop**

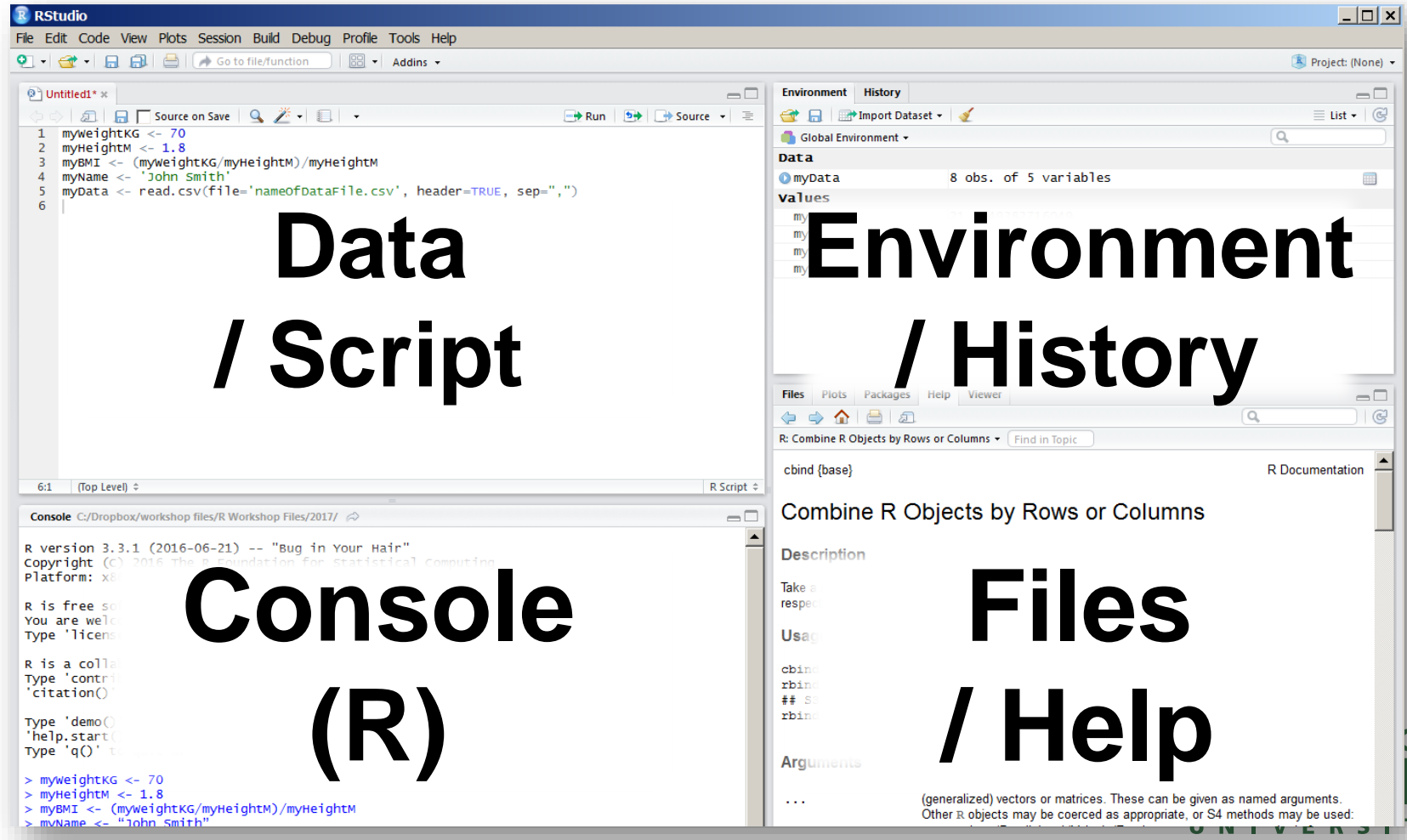https://dsc.gmu.edu/files/Installing-R-RStudio.pdf

**Cloud**

https://posit.cloud

RStudio ➜ Posit

# RStudio Sections

# Benefits of RStudio

- One window to contain everything
- Projects help you contain projects and use a working directory
- Point-and-click for simple tasks (import files, install packages)
- See objects in environment.
- Autocompletion for packages, functions, and objects
- Help documentation at your fingertips
- Easy to connect to a git repository.

# Packages

- Packages = Groups of functions

- Some are built-in to R

- Most are written by researchers
  - Make it easier to do something they want to do
  - Must be installed, just once

- Packages may have functions with the same name!
  - Must tell R which package to use.
  - library(**package**)
  - require(**package**)
  - **package::**function()

**System Library**

| | | |
|---|---|---|
| ☑ | base | The R Base Package |
| ☐ | boot | Bootstrap Functions (Originally by Angelo Canty for S) |
| ☐ | class | Functions for Classification |
| ☐ | cluster | "Finding Groups in Data": Cluster Analysis Extended Rousseeuw et al. |
| ☐ | codetools | Code Analysis Tools for R |
| ☐ | compiler | The R Compiler Package |
| ☑ | datasets | The R Datasets Package |
| ☐ | foreign | Read Data Stored by 'Minitab', 'S', 'SAS', 'SPSS', 'Stata', 'Systat', 'Weka', 'dBase', ... |
| ☑ | graphics | The R Graphics Package |

**User Library**

| | | |
|---|---|---|
| ☐ | abind | Combine Multidimensional Arrays |
| ☐ | afex | Analysis of Factorial Experiments |
| ☐ | anytime | Anything to 'POSIXct' or 'Date' Converter |
| ☐ | askpass | Safe Password Entry for R, Git, and SSH |
| ☐ | assertthat | Easy Pre and Post Assertions |
| ☐ | backports | Reimplementations of Functions Introduced Since R-3.0.0 |
| ☐ | base64enc | Tools for base64 encoding |
| ☐ | BayesFactor | Computation of Bayes Factors for |

# R's Three Areas:

Data Management

Statistical Modeling

Visualization

GEORGE MASON UNIVERSITY

# R for Data Management

- Import Data
- Tidy Data
- Transform Data

# Open Refine - https://openrefine.org/

# Importing Data

- CSV / Delimited
- Excel
- SPSS, Stata, SAS
- json, xml
- packages

*see other sections*

# R for Data Management

| Base R | Tidyverse | Data Table |
|---|---|---|
| data.frame | tibble | data.table |
| | | |
| df$word | %>% | dt[.., .., ..] |
| df[…] | Hadley Wickham | |
| | RStudio/Posit | |

# Tidyverse

**https://www.tidyverse.org/**

**dplyr**
    Select
    Filter
    Mutate
    Summarize
    Arrange

# The pipe

Tidyverse functions take the data as the 1$^{st}$ argument.
This makes it possible to use the pipe.

*function*(**data**, *arguments*)

**data %>%** *function*(*arguments*)

**data >|** *function*(*arguments*)

# More Tidyverse Packages

- For working with specific data types

# R for Statistics

- Formula Notation
- Model Objects
- Output Formatting

GEORGE MASON UNIVERSITY

# jamovi – https://jamovi.org

# General Linear Models

- Regression /ANOVA tools are part of base R.

- Uses formula notation (mostly the Wilkinson Notation)

- The name of the data table object is given as an argument.

**myanova** **<- aov(** fare ~ class + gender**,** **data = titanic** **)**

**mylogistic** **<- glm(** survived ~ class * gender, **data = titanic** ,

**family** = binomial **)**

| ~ | predicted from | : | interaction |
|---|----------------|---|-------------|
| + | include ("and") | * | factorial |

GEORGE MASON UNIVERSITY

# Packages for Statistics

- summarytools – Descriptive Statistics

- afex – ANOVA (also mixed models with lme4)

- emmeans – Postestimation tests

- lme4 – Mixed Models

- lavaan – Structural Equation Modeling

# Output

- Functions create an analysis object, not useful by itself
- Use other functions to display tables, plots, etc
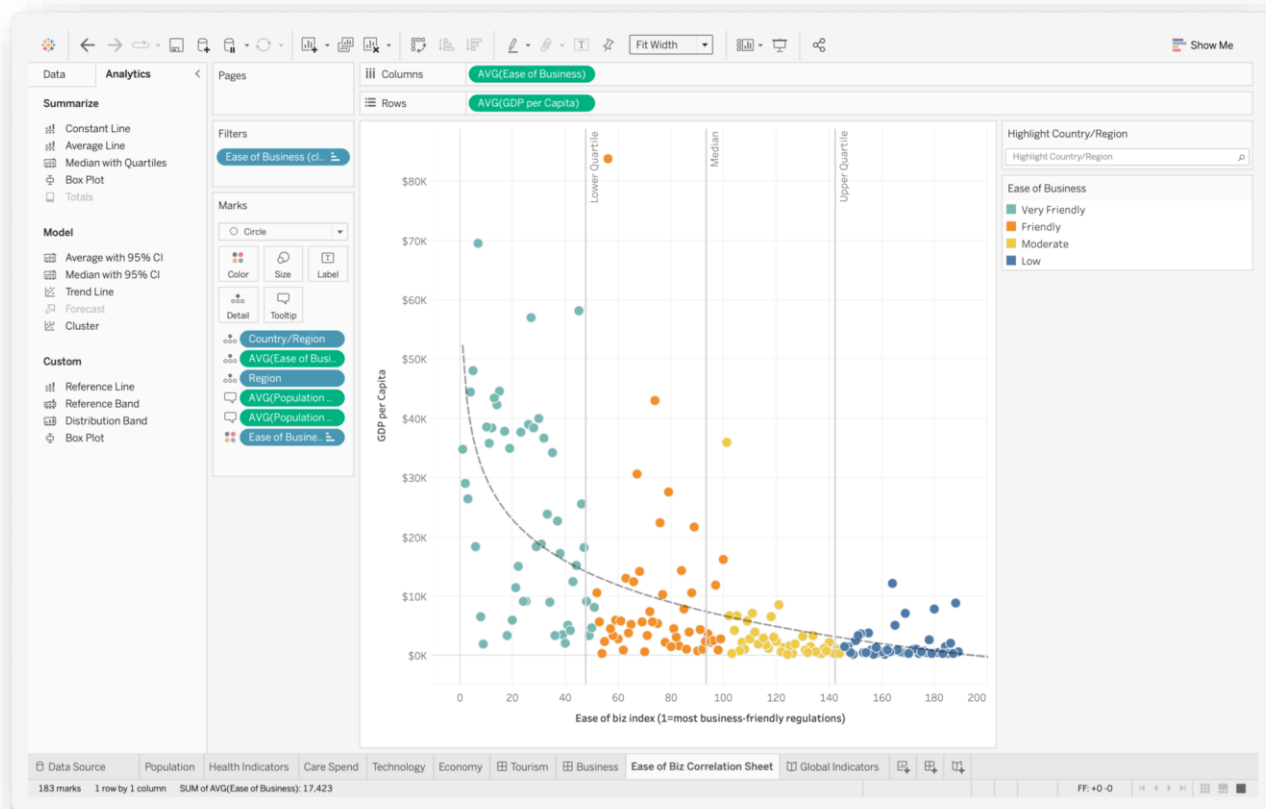- e.g., summary() function for lm / glm / aov



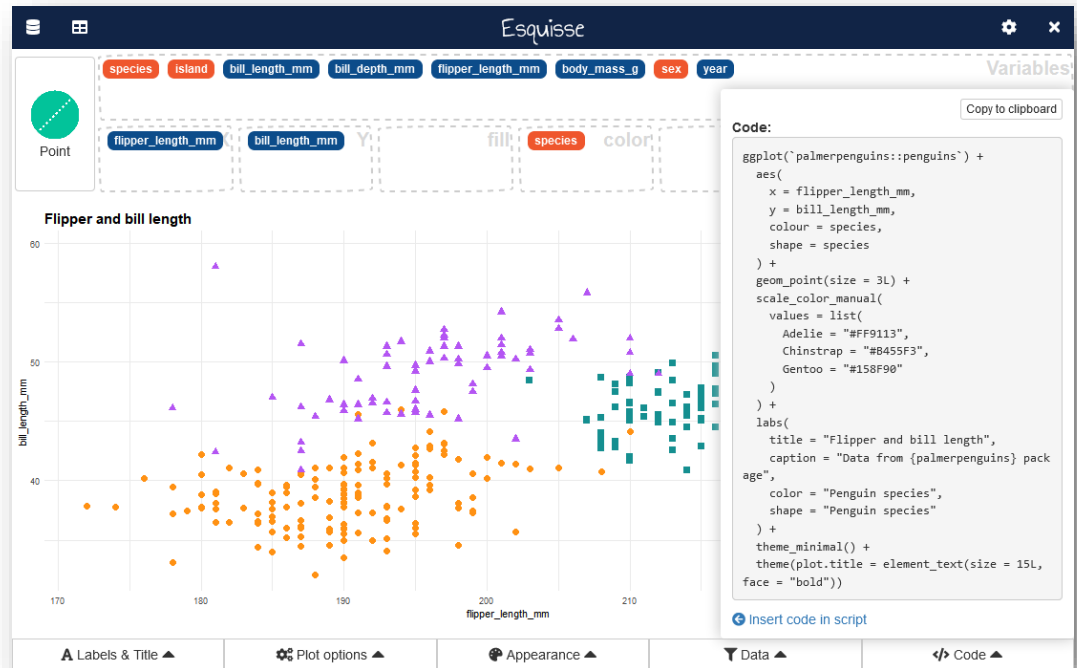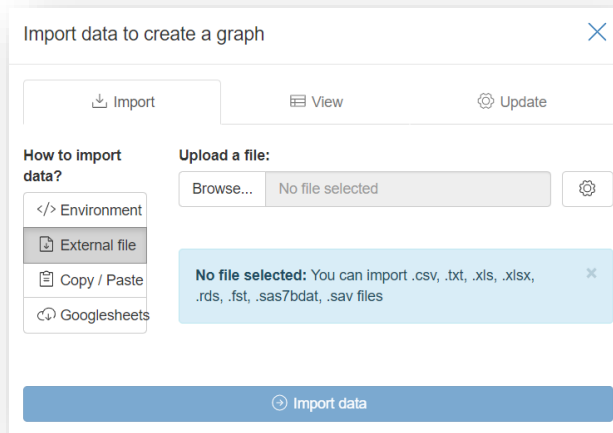stargazer, gt / gtsummary, kable

# R for Graphing

# Tableau

Free for Academic Use: https://www.tableau.com/academic/students
Free for Any Use: https://public.tableau.com (can only save online)

# esquisse - https://dreamrs.github.io/esquisse/

```
install.packages("esquisse")
esquisse::esquisser()
```

# ggplot (ggplot2)

- https://ggplot2.tidyverse.org
- https://r-graph-gallery.com/ggplot2-package.html
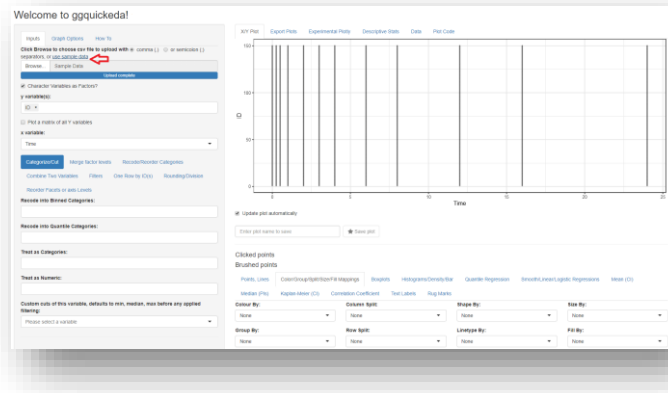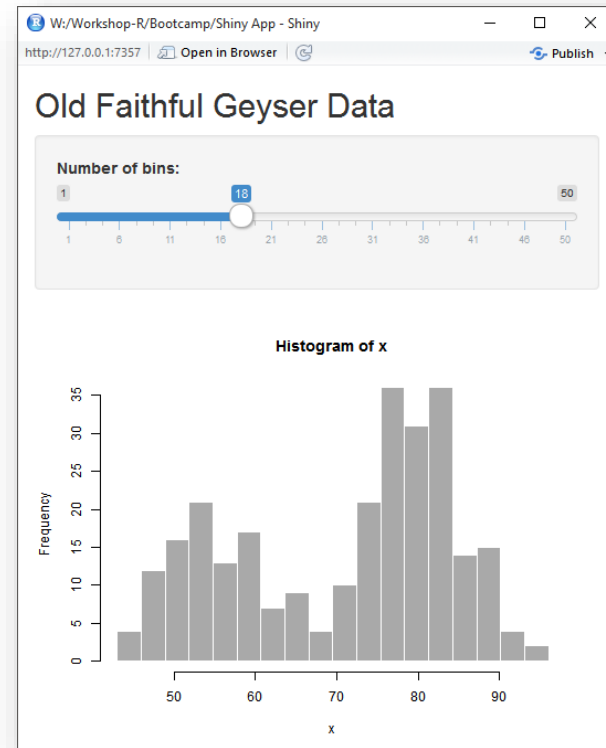
# Easier Graphs

- [ggquickeda](#) – another Point and Click Interface



- [GGally](#) - Good for bivariate graphing with functions like ggpairs()

# Interactive Graphics (Shiny)

# Notebooks

- Quarto Document

- Jupyter Notebook

# Reading Tutorials

\#      = A Comment

\#>    = Commented output (what you should see)

            Does NOT save output in an object

df <- tibble(…)    = Creates a dataset

<u>word</u>  =  a link, click to learn more

```
# By default, mutate() keeps all columns from the input data.
df <- tibble(x = 1, y = 2, a = "a", b = "b")
df %>% mutate(z = x + y, .keep = "all") # the default
#> # A tibble: 1 × 5
#>         x     y a       b                z
#>     <dbl> <dbl> <chr>  <chr>        <dbl>
#> 1      1     2 a       b                3
df %>% mutate(z = x + y, .keep = "used")
#> # A tibble: 1 × 3
```

# Finding Tutorials

To get started

- Look for ones that do what you need to do, preferably from someone in your field.

- Follow along!

- primers, learnr, swirl

- Many free online books

- Audit MOOCs

- Paid sites offer some free interactive tutorials

GEORGE MASON UNIVERSITY

# Finding Tutorials

For doing specific tasks

- Figure out what package you will want to use
  - Colleagues
  - In your field
- Look at "vignettes" and documentation for the package
- See recommended sites

GEORGE MASON UNIVERSITY

# Ignore tutorials that

- Don't use your data management scheme
  - e.g., if it uses the functions apply(), sapply(), tapply(), that is base R

- Look complicated

- Are more than 2 years old (generally)

# Data for Follow-Along Tutorials

Use data loaded with base R or another package

Create data using functions or by hand

Provide a data file to download and read in

# Teaching with R

- Posit Cloud
  - Free plan with limited sharing may be sufficient
  - Instructor plan with unlimited sharing is $15 + $5/month/student

- DataCamp
  - Free for students in a class, supports assignments
  - https://www.datacamp.com/groups/classrooms

- OER
  - Many OER textbooks (most online), more being added
  - Including https://bookdown.org
  - https://education.rstudio.com/teach/materials/