

PSYC 259:

Principles of Data Science

Week 2: File Organization and
Workflow

Today

1. Project structure principles
 - a. File/folder organization
 - b. Version control
2. Advice: How to get programming help
3. BREAK
4. R language basics and importing tutorial
5. Getting started on homework

Project Structure Principles

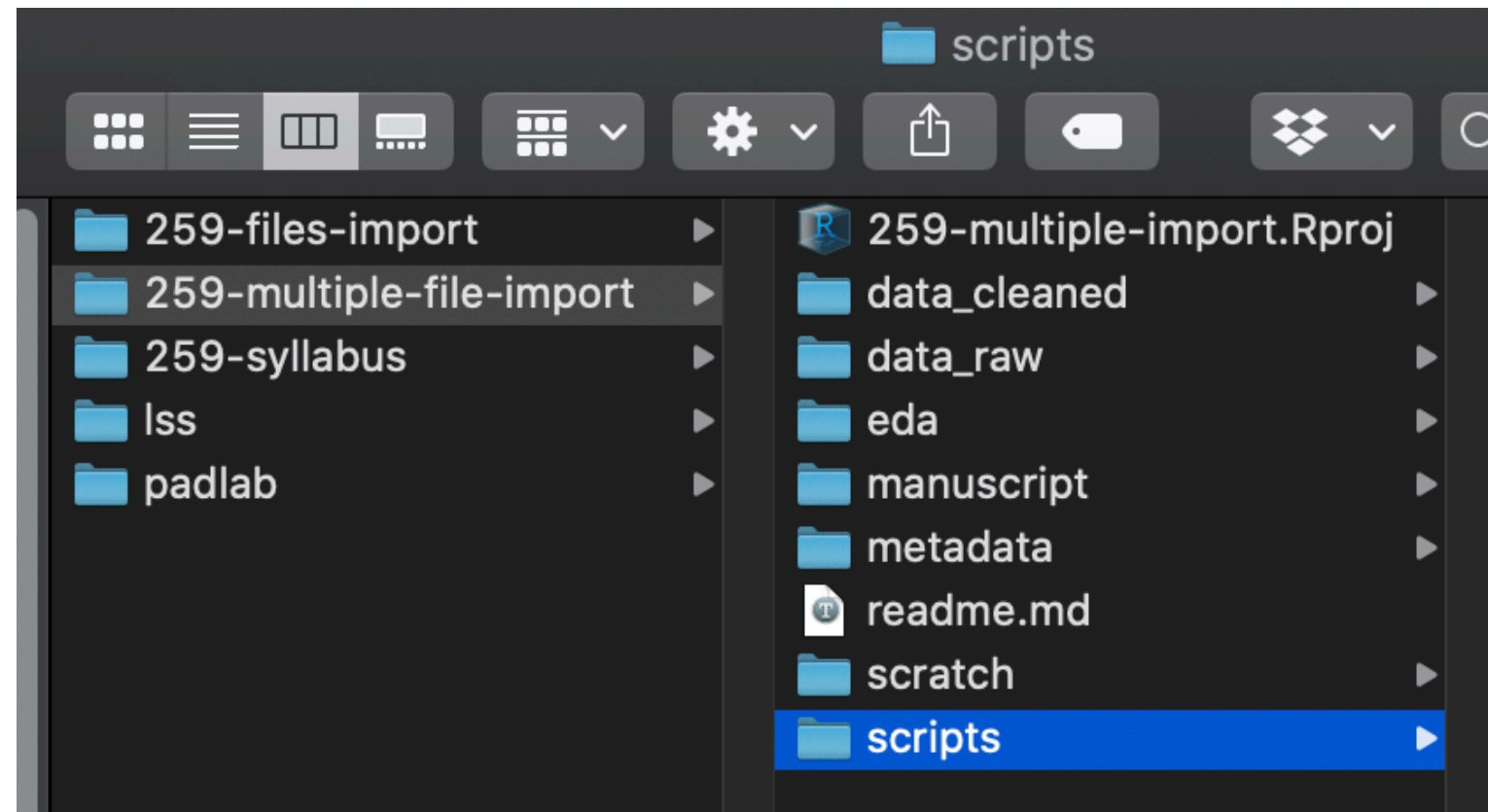
Two principles of project workflow

1. Folder organization creates rules and defines a workflow; establishes a location from which to build relative file paths
2. Version control tracks file history without duplication/clutter; allows for collaboration/derivation/experimentation

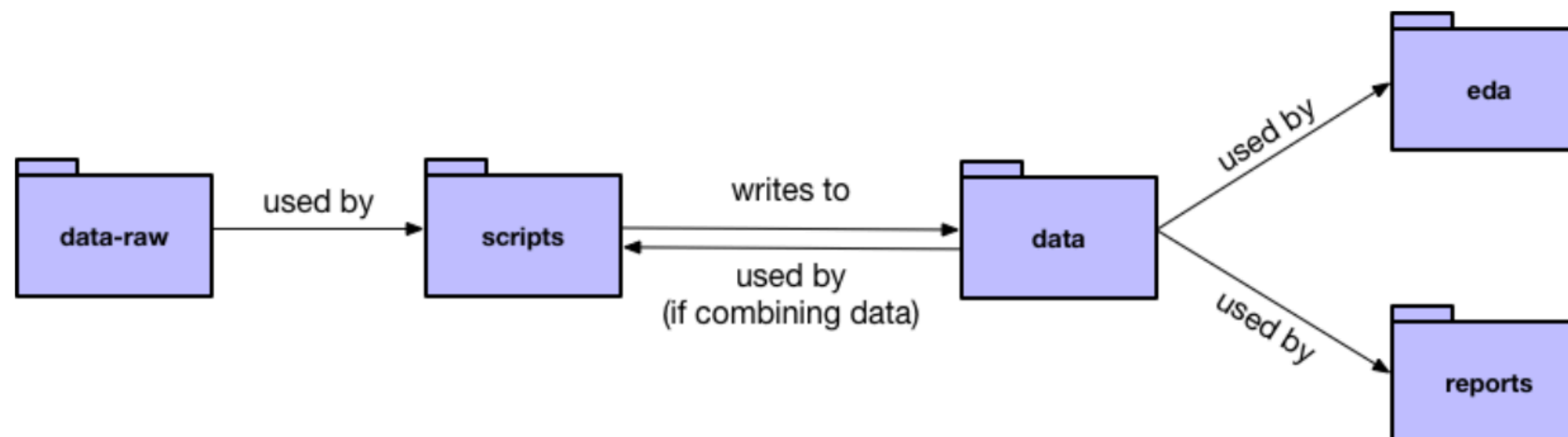
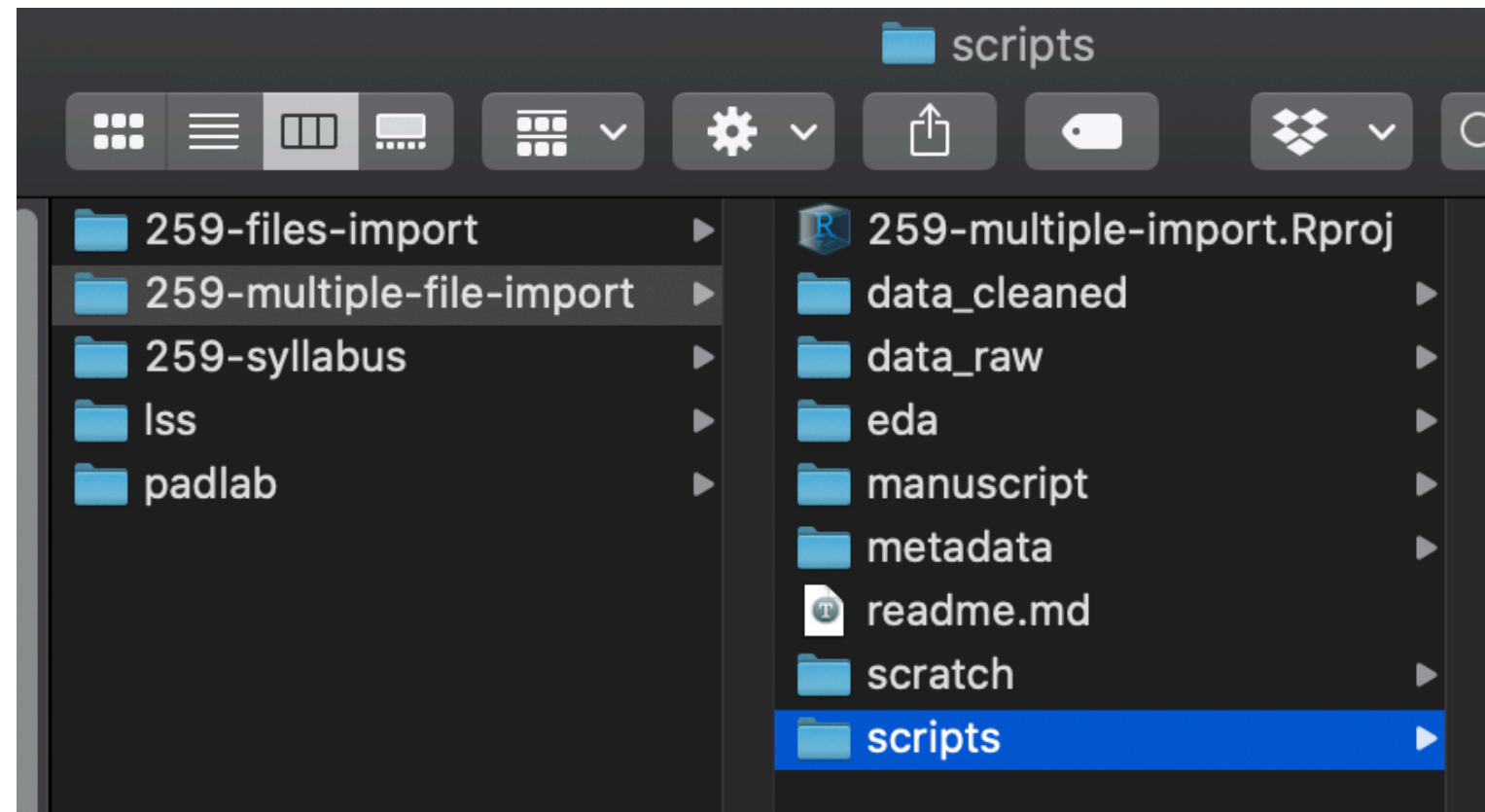
What principles should
guide project structure?

#1 File/folder organization

Folders should organize similar file **types** (w/in a root project folder)

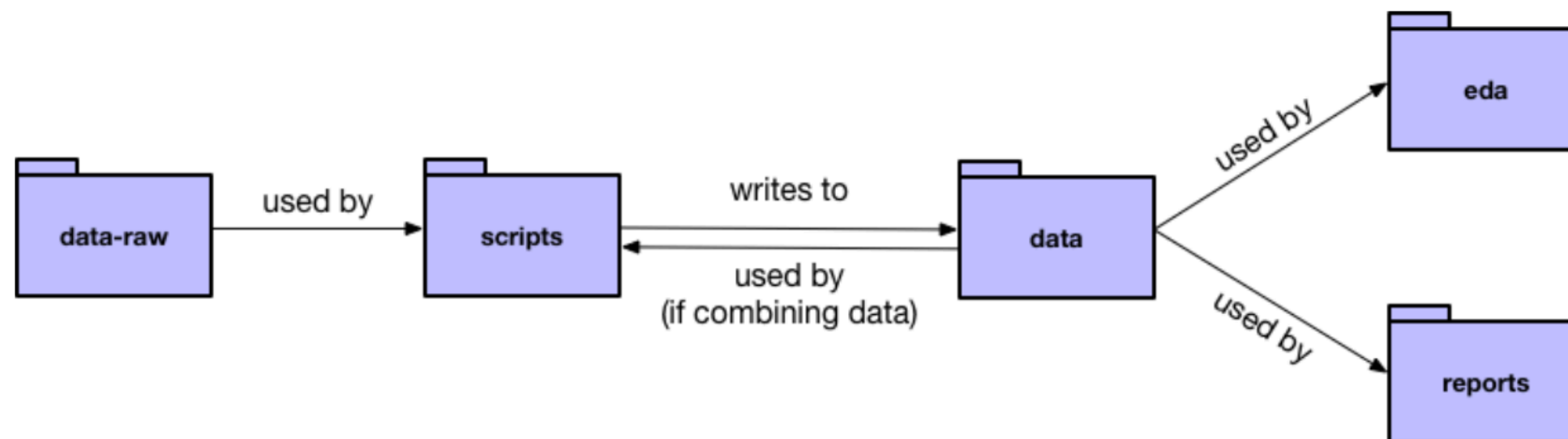


Folders should organize similar file **types** (w/in a root project folder)



Avoid having more nodes in this chart than are necessary

- Every intermediate step incurs a cost of maintenance
 - Think about what has to be re-run if data-raw changes
- Multiple endpoints don't incur cost



What is an R Studio Project?

- A project folder (directory) with some files that keep tabs on your command history (.RProj)
 - Open the RProj file to open the project
 - Can turn existing folders into projects
- What is real (e.g., persists after shutting down R)?
 - Not real: objects/data frames in your workspace
 - Real: data files and the scripts (.R files) used to work with them

Accessing files from your working directory

- Your RStudio Project folder is your default *working directory*
 - working directory: where R can find files
 - `getwd()` probably gives you something ugly like `/Users/johnfranchak/Documents/GitHub/project_folder`
- *Absolute file paths* like that should be avoided at all costs!
 - Absolute file paths don't transfer between computers with different user names or different operating systems
 - Bad for reproducibility, extensibility, and sharing
- But don't you need them to get into all of those directories you just told me I need to use?

Relative files paths 🦵 🦵 🦵

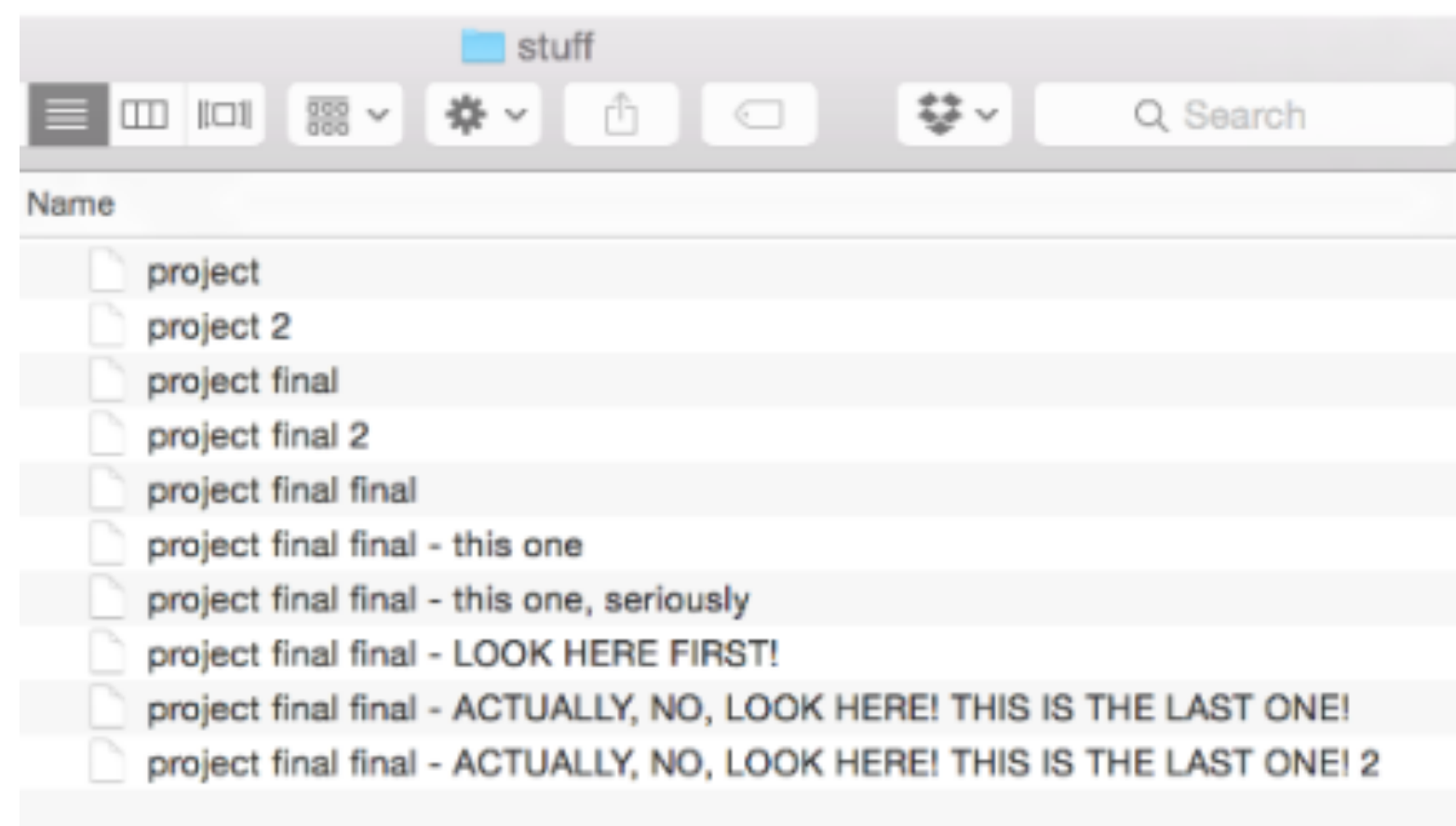
- Just tell R what to look for relative to your project (working) directory!
- append subfolder name to filename
 - "folder1/filename"
- *here* package detects the project directory and composes filenames from/to any folder
 - `here("folder1", "folder2", "filename")`

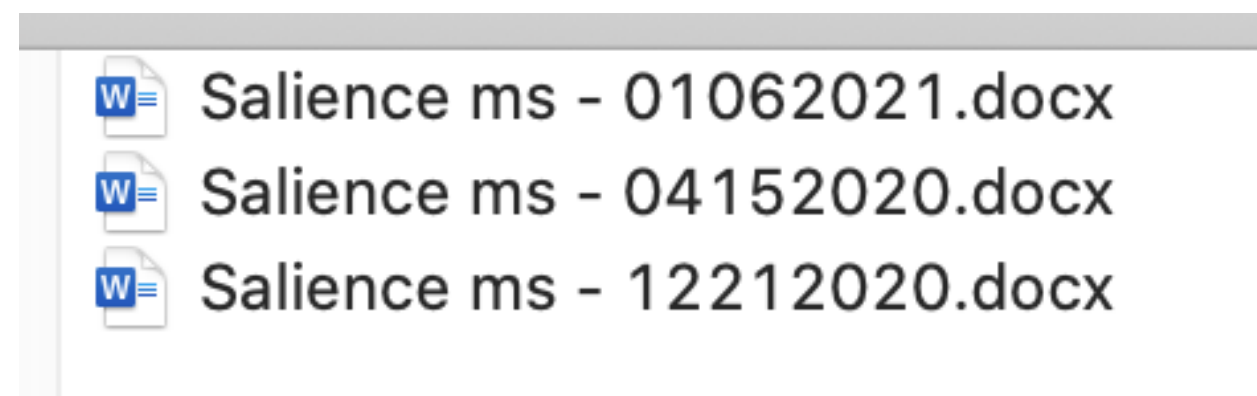
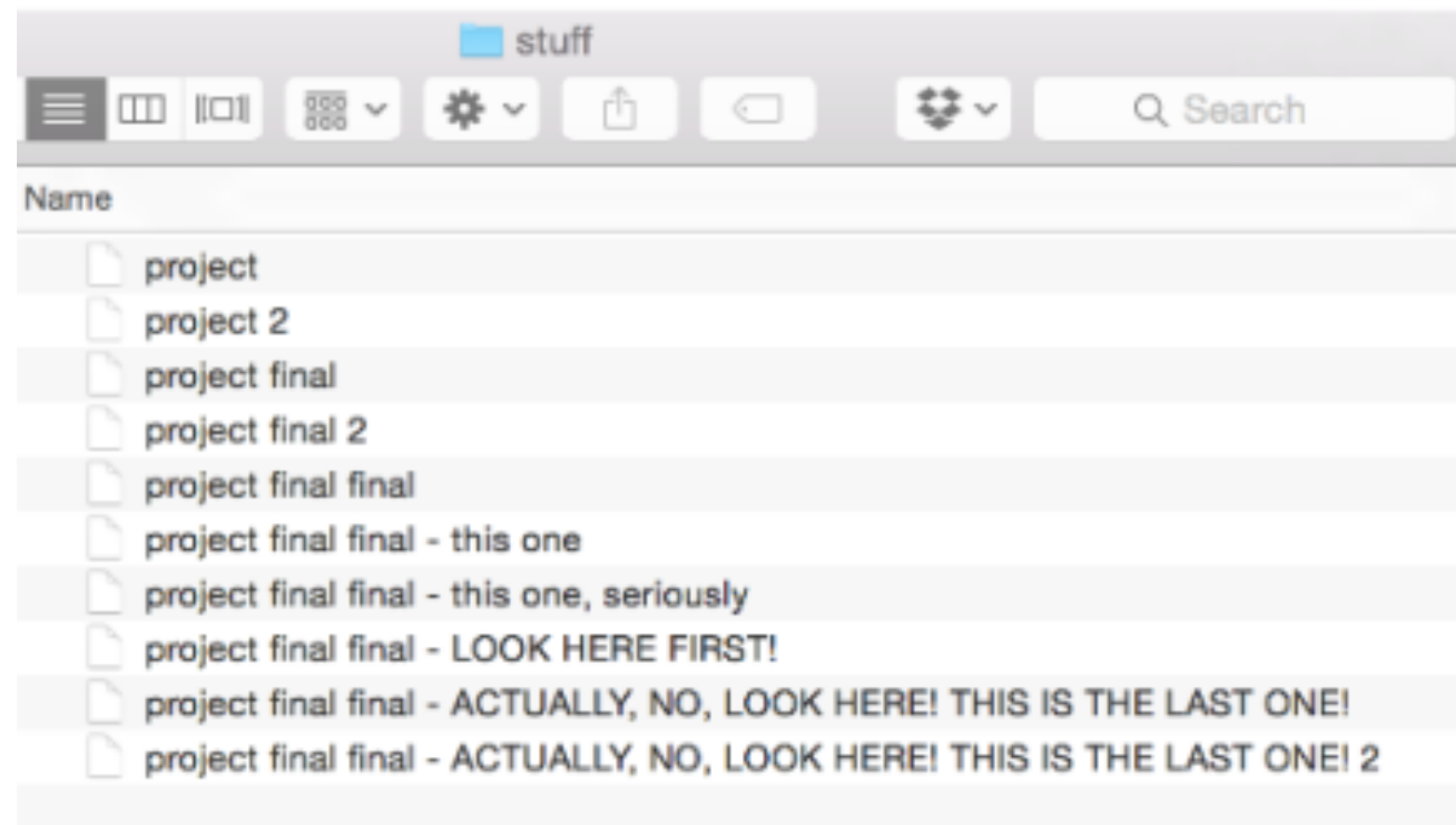
Other considerations

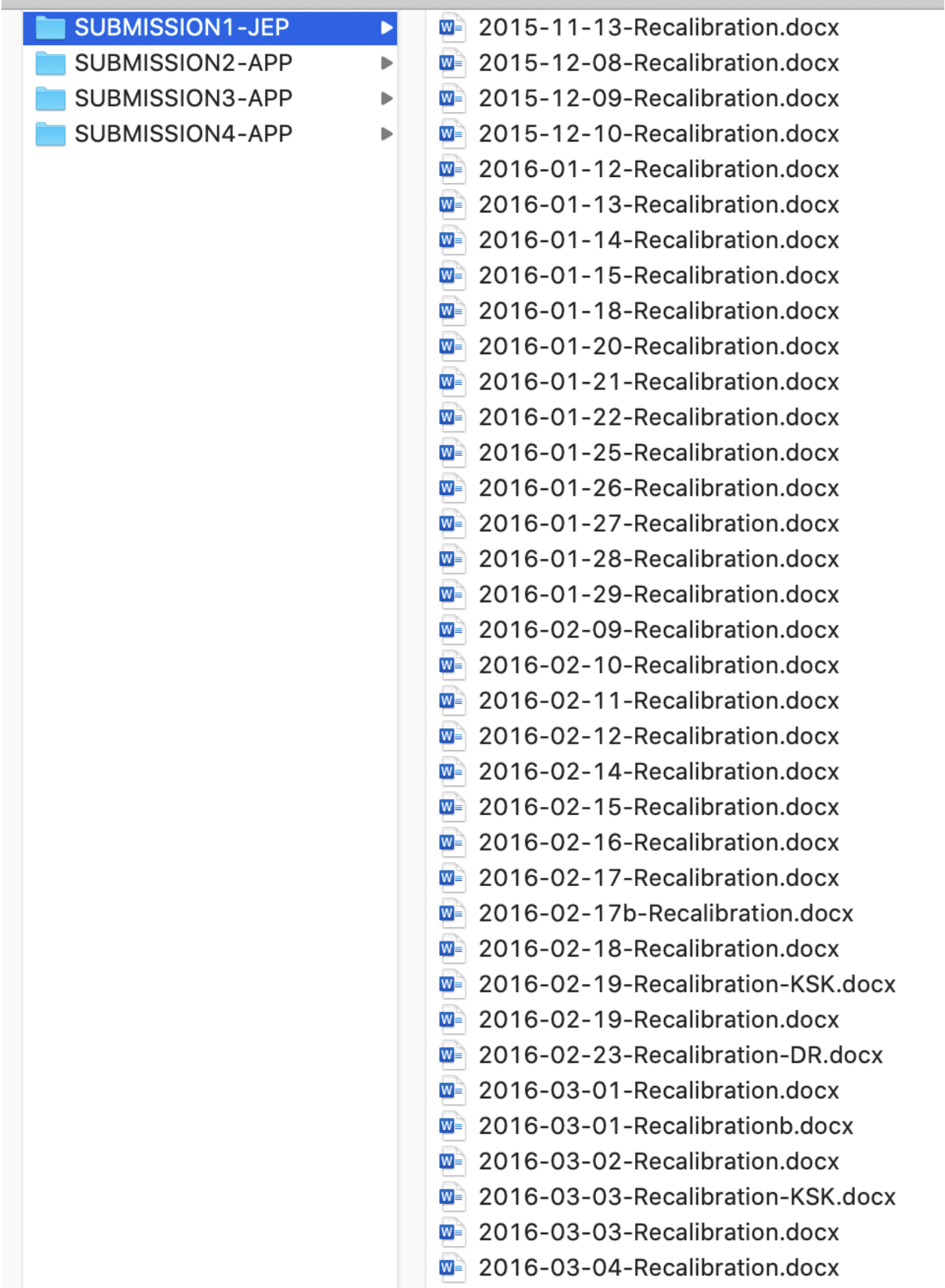
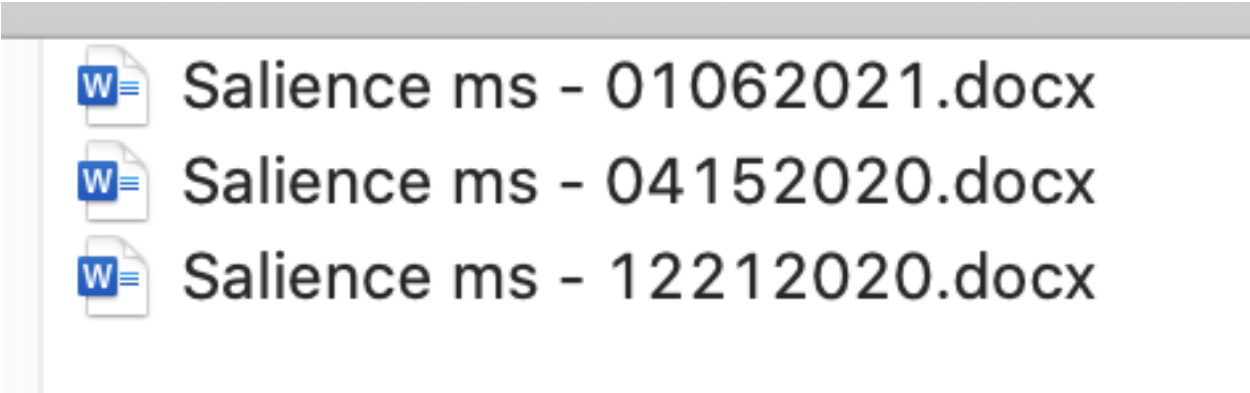
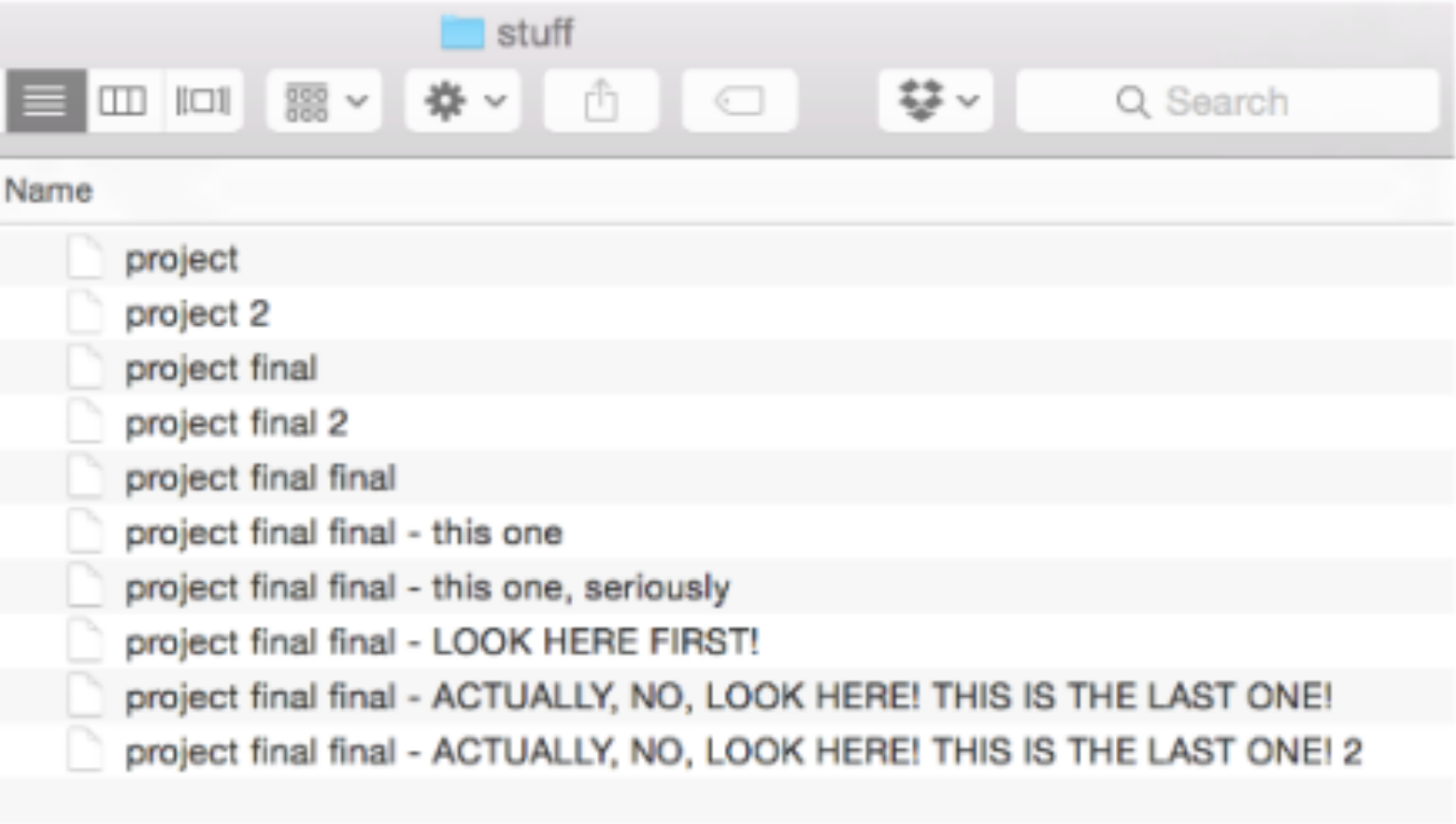
- Not everything can be automated, but you can try to limit human data entry to a single master table
 - Track notes about sessions, inclusion/exclusion info
 - Keep as part of project metadata, and use it to direct your scripts to pull the 'right' data
- Not every project can be contained in a local directory on a single computer (or on github)
 - Large datasets might need other solutions
 - "data_raw" might need to be "data_slightly_cooked"

What principles should
guide project structure?

#2 Version Control







Why do we do this to ourselves?

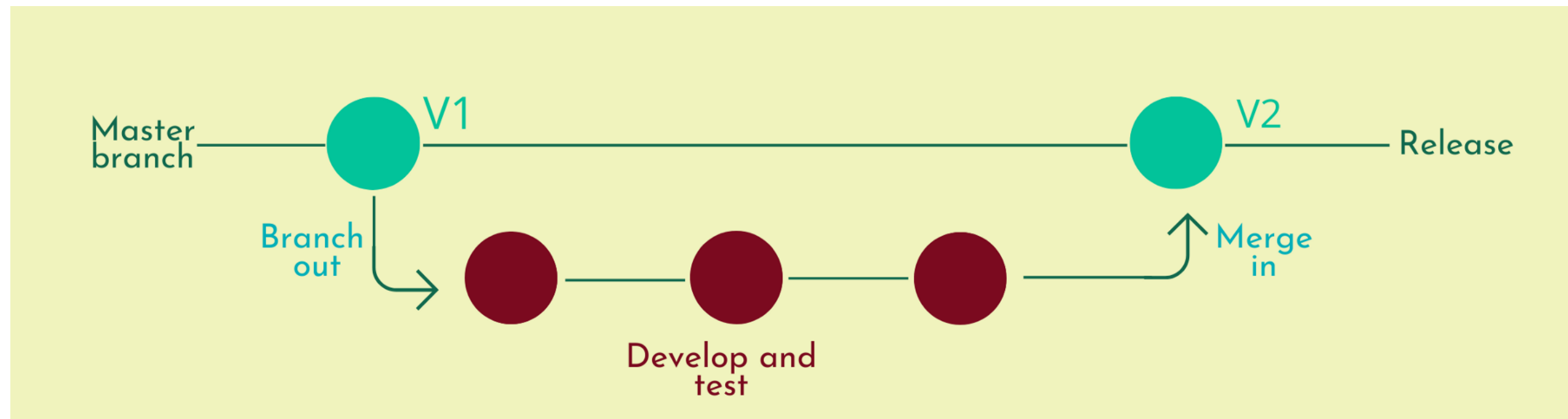
- Want to preserve the document history
 - Easier to take risks, delete things, etc. if you know you can go back (especially for writing code)
 - For collaboration, want to know who made changes
- But it's messy, inefficient, and easy to break
 - Cluttered file folders (imagine if you wanted to track every change to every type of file)
 - Copy -> paste -> rename is tedious
 - No record of **what** the changes were

The solution: Use version control software (e.g., git, svn)

- Define a **repository** (your project folder) and tell git which files to **track** vs **ignore**
- Place your repository on a cloud hub (e.g., GitHub)
- Git tracks what/when changes were made, and you tell git whether to **commit** those changes
- You can **push** those changes to the central hub to store them (now they're saved)
- Other users/computers can **pull** those changes from the hub to keep their local copy up-to-date

More advanced git features

- Create branches (alternate timelines) to develop and test new functions, then merge those changes back to the master branch when they're ready



Other considerations

- Easiest to get started working on a solo project
- Collaborative version control can get messy if people forget to commit/push their changes or to pull the most recent version from the master repo
- git wants your folders out of dropbox or other cloud synced drives, so you have to commit to using it
- Lots of ways to use git (Github website, Github app, RStudio git window, command line)
- Public vs. private repos
- Flat files vs. binary files

Essential git skills to master

- Forking repos - creating a copy of a repo that points to the original [DONE]
- Cloning a repo - pulling a version from the cloud onto your local system [DONE]
- Pulling changes from the cloud
- Committing and pushing changes to the cloud
- Discarding unwanted changes
- Creating a new repo from an existing local directory

259-langbasics-importing Public

forked from [psych-259-data-science-2022/259-langbasics-importing](#)

Watch 0

Fork 1

Star 0

Fork your copy of 2025-PSYC-259/259-langbasics-importing

master

1 Branch 0 Tags

Go to file

Add file

Code

This branch is 3 commits ahead of [psych-259-data-science-2022/259-langbasics-importing:master](#).

JohnFranchak fixed demos ✓	88fdf1c · last week	🕒 33 Commits
data_cleaned	made a standalone version of the script examples	4 years ago
data_raw	made a standalone version of the script examples	4 years ago
docs	fixed demos	last week
.gitignore	created flipbook	4 years ago
259-langbasics-importing.Rproj	created flipbook	4 years ago
README.md	2025	3 weeks ago
langbasics-importing.R	fixed demos	last week
xaringan-themer.css	2025	3 weeks ago

About

Introduction to R data importing

[2025-PSYC-259.github.io/259-langba...](#)

Readme

Activity

Custom properties

0 stars

0 watching

1 fork

Report repository

Releases

No releases published

Packages

No packages published

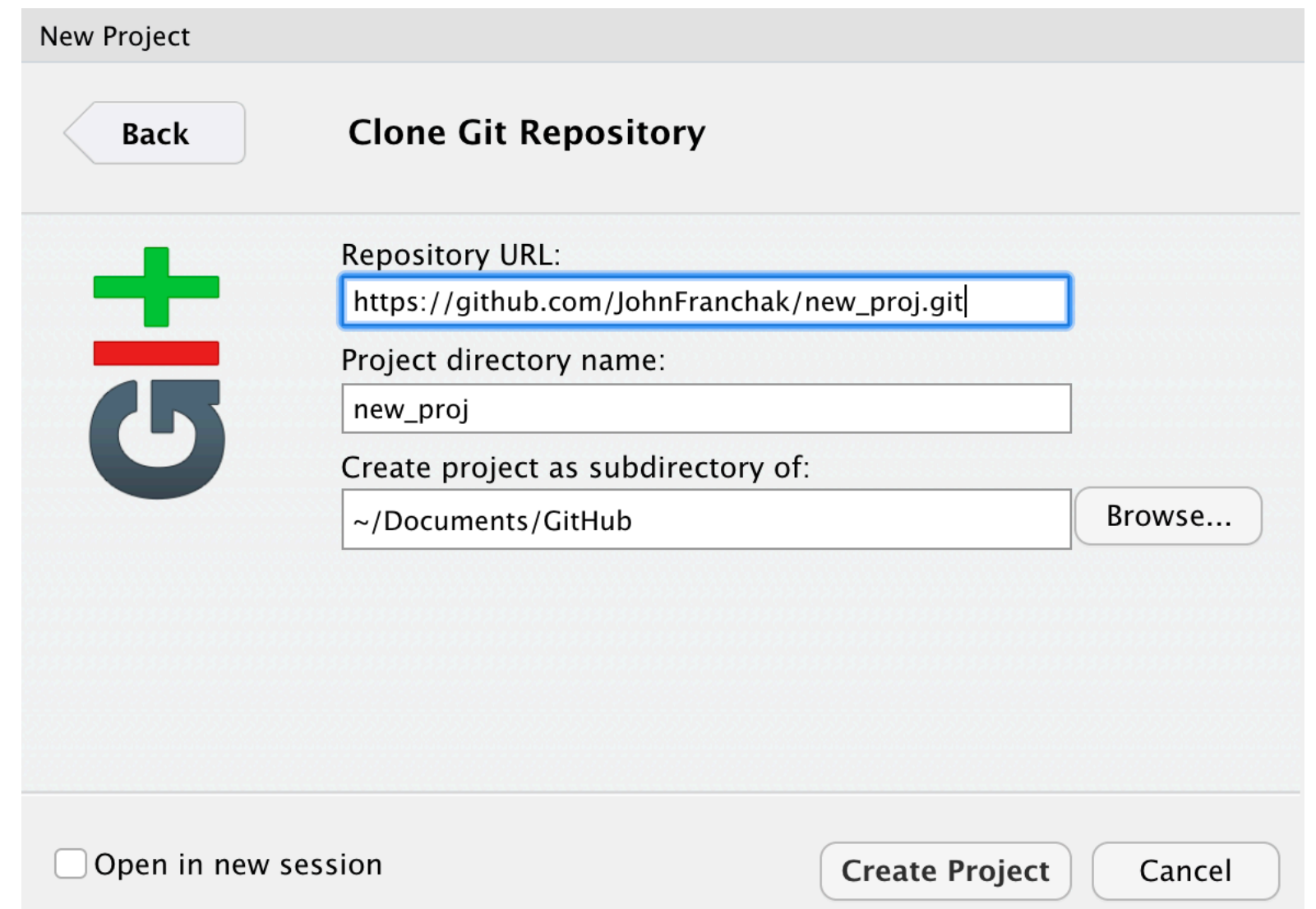
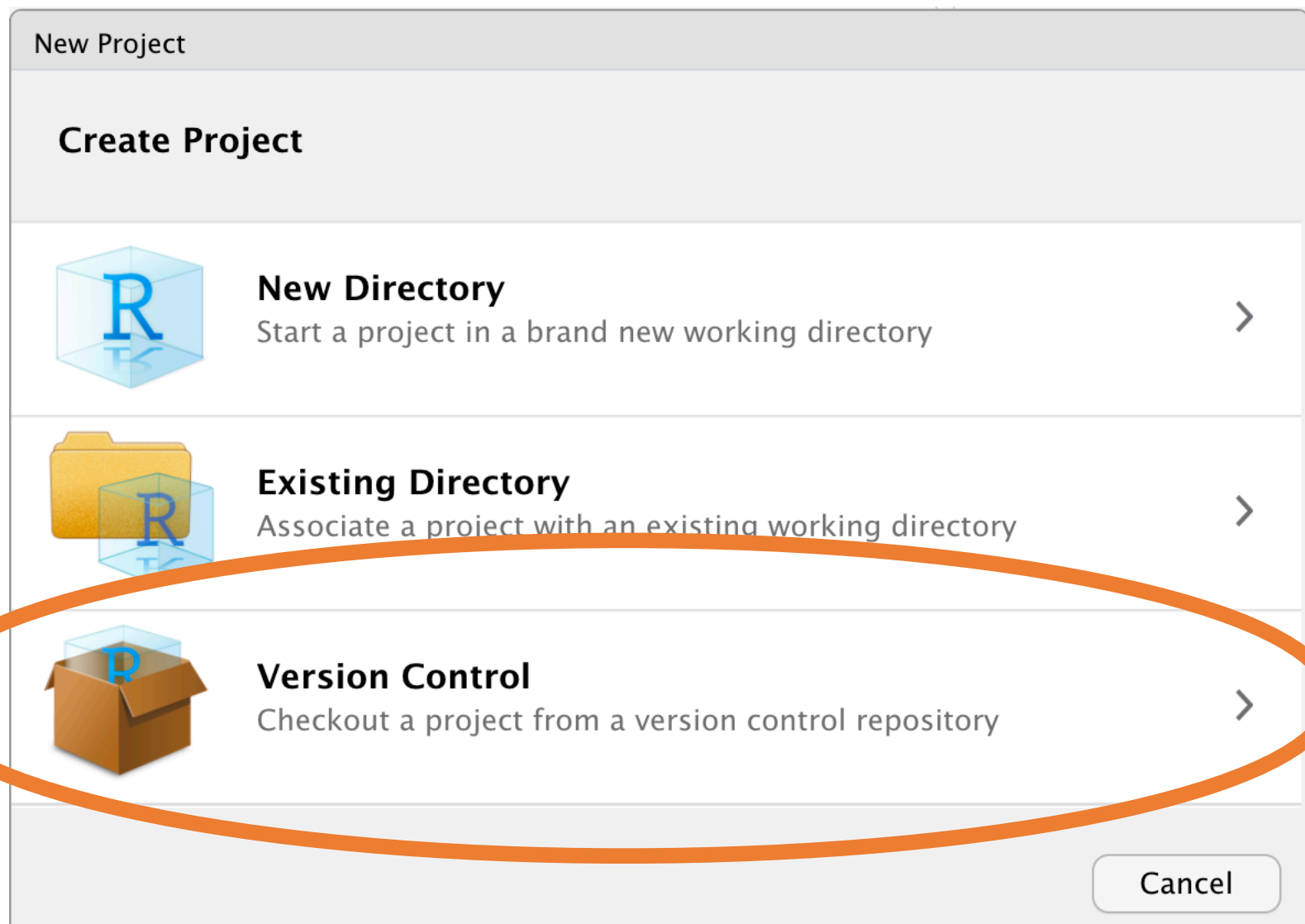
Fork a repository to make your own copy of it

- Log in to github.com with your username
- Go to <https://github.com/2025-PSYC-259/259-langbasics-importing>
- Now it won't be PSYC-259-Data-Science's anymore, it will be yours!
- Let's get it off the Github website and onto your local computer by **cloning it**

Disclaimer!

- We are not going over how to work with the Github Desktop app here. We will just be working with GitHub online and using RStudio to connect to GitHub

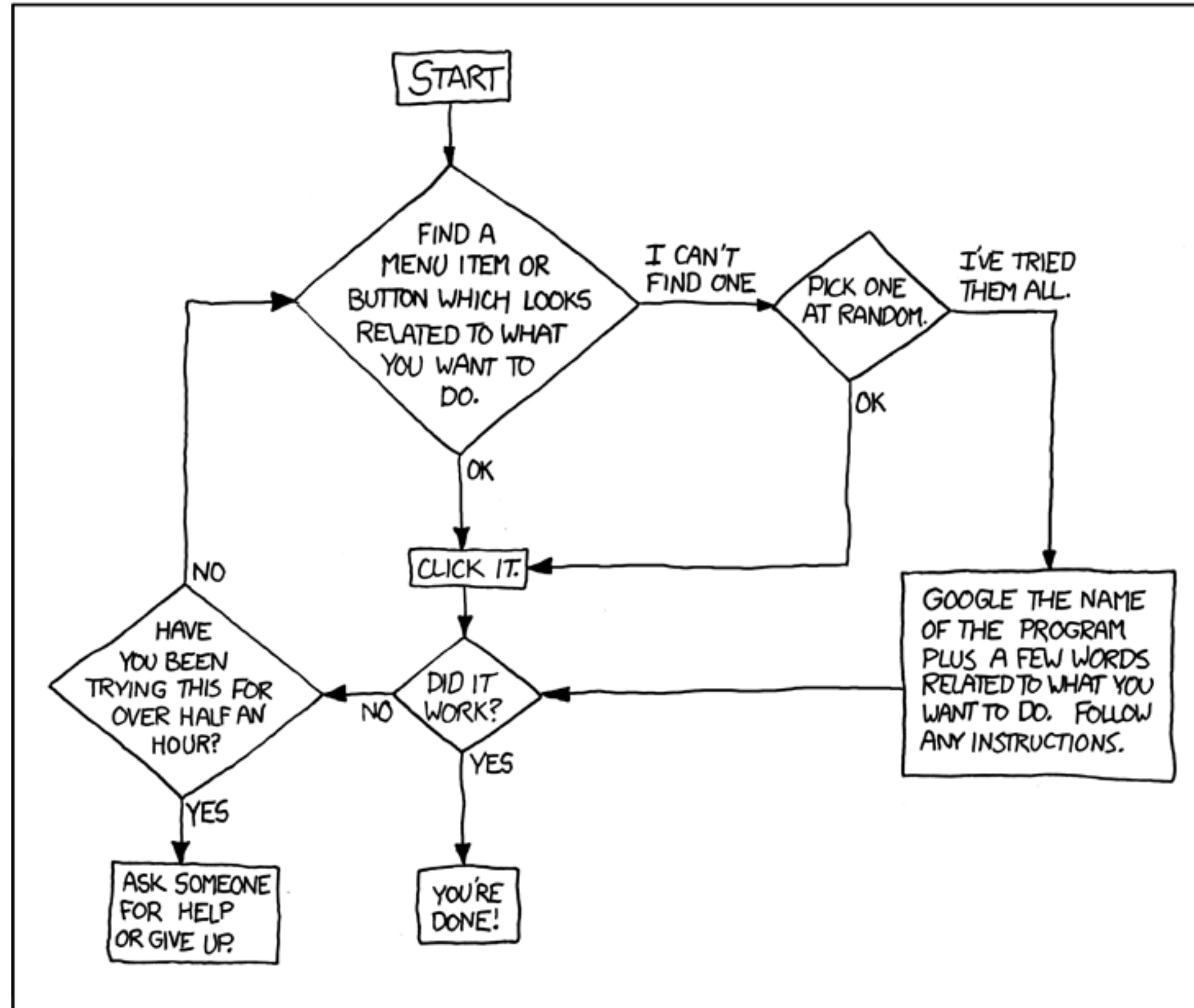
Clone a project directly into RStudio from Github



General advice:
How to get help

DEAR VARIOUS PARENTS, GRANDPARENTS, CO-WORKERS,
AND OTHER "NOT COMPUTER PEOPLE."

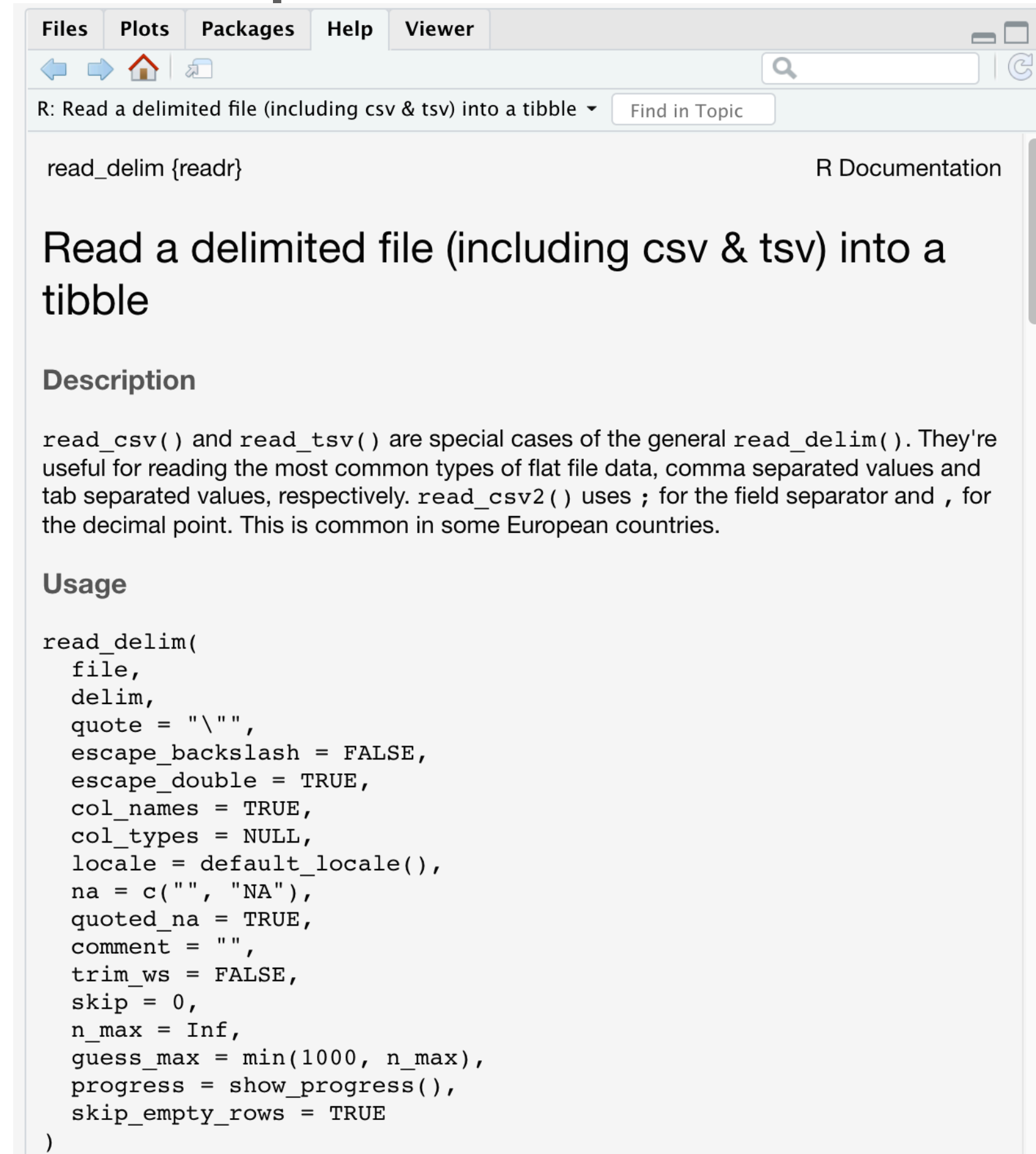
WE DON'T MAGICALLY KNOW HOW TO DO EVERYTHING IN EVERY
PROGRAM. WHEN WE HELP YOU, WE'RE USUALLY JUST DOING THIS:



PLEASE PRINT THIS FLOWCHART OUT AND TAPE IT NEAR YOUR SCREEN.
CONGRATULATIONS; YOU'RE NOW THE LOCAL COMPUTER EXPERT!

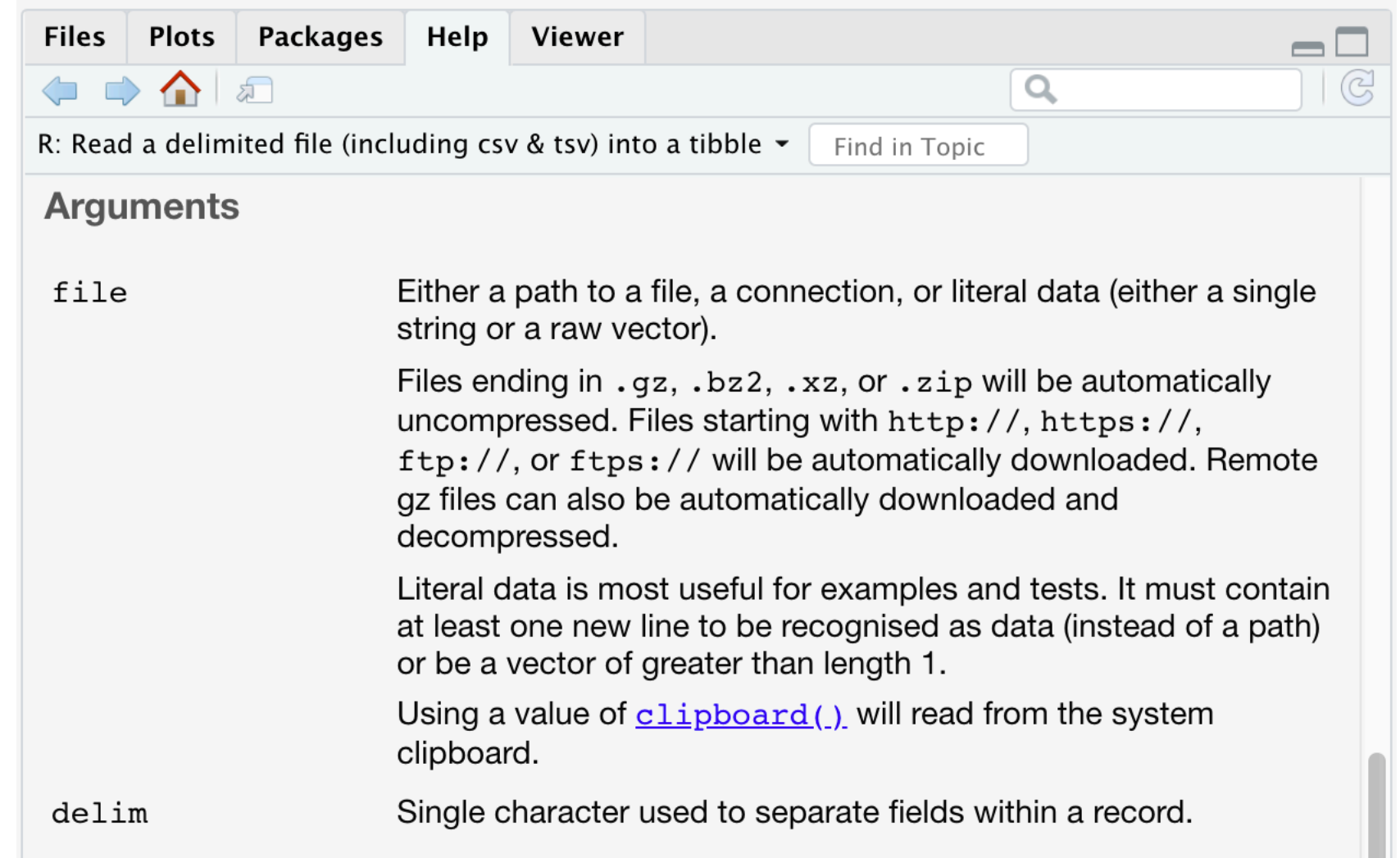
Where to search for help

- R Documentation
 - ?function brings up the documentation for a function (assuming the package is loaded)



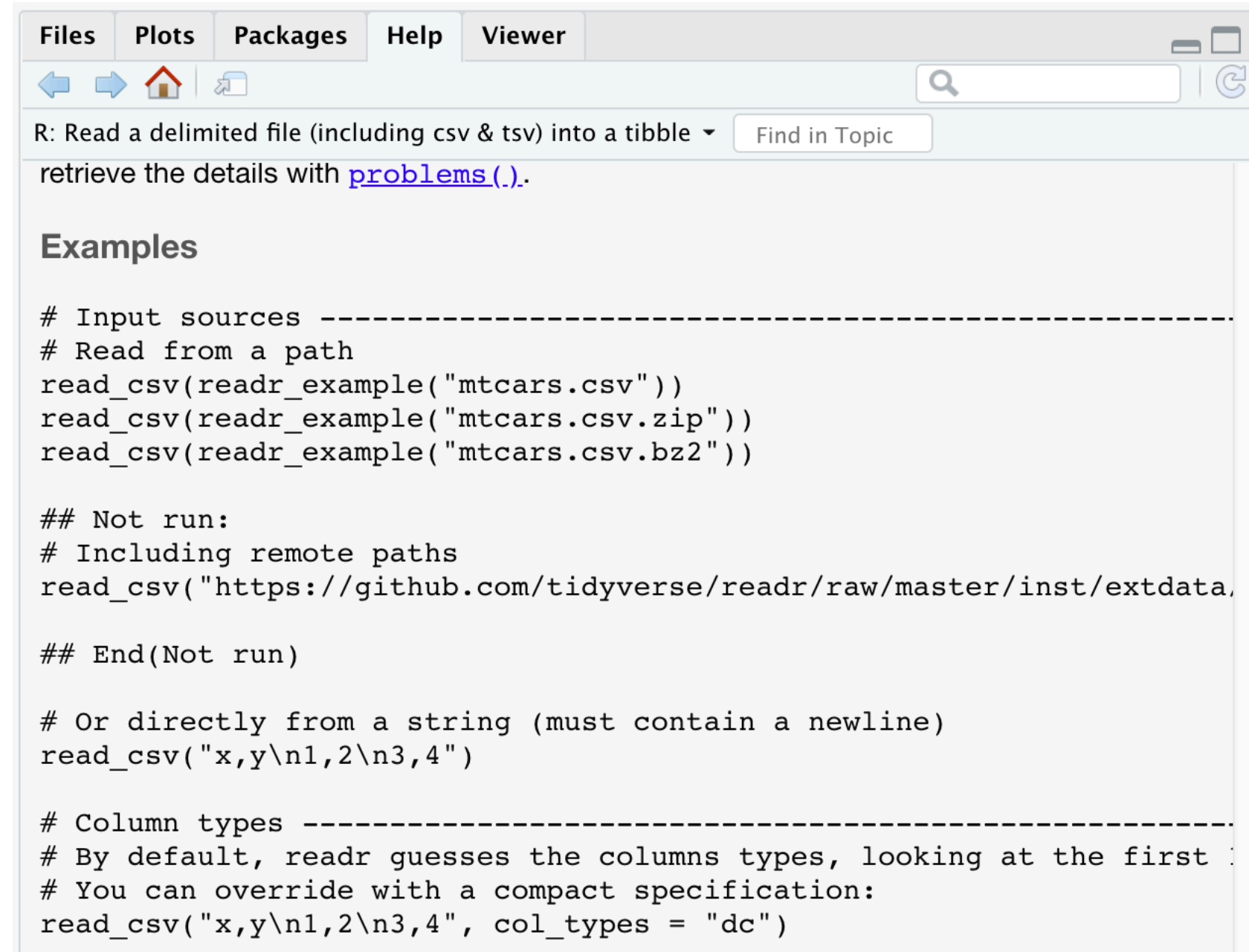
Where to search for help

- R Documentation
 - ?function brings up the documentation for a function (assuming the package is loaded)
 - Scroll down to see the “arguments” definitions



Where to search for help

- R Documentation
 - ?function brings up the documentation for a function (assuming the package is loaded)
 - Scroll down to see the “arguments” definitions
 - Scroll down more to see examples



The screenshot shows the RStudio Help window for the `read_csv` function. The window has a menu bar with 'Files', 'Plots', 'Packages', 'Help', and 'Viewer'. Below the menu bar is a toolbar with navigation icons and a search bar. The main content area displays the documentation for 'R: Read a delimited file (including csv & tsv) into a tibble'. It includes a 'Find in Topic' button and a search bar. The text describes how to retrieve details with `problems()`. Under the 'Examples' section, there are several code snippets demonstrating how to use `read_csv` with different input sources, including local paths, remote URLs, and strings. The examples are separated by dashed lines.

```
# Input sources -----
# Read from a path
read_csv(readr_example("mtcars.csv"))
read_csv(readr_example("mtcars.csv.zip"))
read_csv(readr_example("mtcars.csv.bz2"))

## Not run:
# Including remote paths
read_csv("https://github.com/tidyverse/readr/raw/master/inst/extdata/

## End(Not run)

# Or directly from a string (must contain a newline)
read_csv("x,y\n1,2\n3,4")

# Column types -----
# By default, readr guesses the columns types, looking at the first
# You can override with a compact specification:
read_csv("x,y\n1,2\n3,4", col_types = "dc")
```

Where to search for help

- Package vignettes/blogs
 - <https://cran.r-project.org/web/packages/readr/vignettes/readr.html>
- Google
 - Blogs/instructional sites
 - StackExchange (especially for error text)
 - Tidyverse community
 - <https://community.rstudio.com/c/tidyverse/6>
- Take working examples and fiddle!

Tidyverse cheatsheets!

Data Import :: CHEAT SHEET

R's **tidyverse** is built around **tidy data** stored in **tibbles**, which are enhanced data frames.



The front side of this sheet shows how to read text files into R with **readr**.



The reverse side shows how to create tibbles with **tibble** and to layout tidy data with **tidyr**.

OTHER TYPES OF DATA

Try one of the following packages to import other types of files

- **haven** - SPSS, Stata, and SAS files
- **readxl** - excel files (.xls and .xlsx)
- **DBI** - databases
- **jsonlite** - json
- **xml2** - XML
- **httr** - Web APIs
- **rvest** - HTML (Web Scraping)

Save Data

Save **x**, an R object, to **path**, a file path, as:

Comma delimited file

write_csv(x, path, na = "NA", append = FALSE, col_names = !append)

File with arbitrary delimiter

write_delim(x, path, delim = " ", na = "NA", append = FALSE, col_names = !append)

CSV for excel

write_excel_csv(x, path, na = "NA", append = FALSE, col_names = !append)

String to file

write_file(x, path, append = FALSE)

String vector to file, one element per line

write_lines(x, path, na = "NA", append = FALSE)

Object to RDS file

write_rds(x, path, compress = c("none", "gz", "bz2", "xz", ...))

Tab delimited files

write_tsv(x, path, na = "NA", append = FALSE, col_names = !append)

Read Tabular Data - These functions share the common arguments:

read_*(file, col_names = TRUE, col_types = NULL, locale = default_locale(), na = c("", "NA"), quoted_na = TRUE, comment = "", trim_ws = TRUE, skip = 0, n_max = Inf, guess_max = min(1000, n_max), progress = interactive())

a,b,c
1,2,3
4,5,NA

A	B	C
1	2	3
4	5	NA

Comma Delimited Files

read_csv("file.csv")

To make file.csv run:

write_file(x = "a,b,c\n1,2,3\n4,5,NA", path = "file.csv")

a;b;c
1;2;3
4;5;NA

A	B	C
1	2	3
4	5	NA

Semi-colon Delimited Files

read_csv2("file2.csv")

write_file(x = "a;b;c\n1;2;3\n4;5;NA", path = "file2.csv")

a|b|c
1|2|3
4|5|NA

A	B	C
1	2	3
4	5	NA

Files with Any Delimiter

read_delim("file.txt", delim = "|")

write_file(x = "a|b|c\n1|2|3\n4|5|NA", path = "file.txt")

a b c
1 2 3
4 5 NA

A	B	C
1	2	3
4	5	NA

Fixed Width Files

read_fwf("file.fwf", col_positions = c(1, 3, 5))

write_file(x = "a b c\n1 2 3\n4 5 NA", path = "file.fwf")

Tab Delimited Files

read_tsv("file.tsv") Also **read_table**()

write_file(x = "a\tb\tc\n1\t2\t3\n4\t5\tNA", path = "file.tsv")

USEFUL ARGUMENTS

a,b,c
1,2,3
4,5,NA

Example file

write_file("a,b,c\n1,2,3\n4,5,NA", "file.csv")
f <- "file.csv"

1	2	3
4	5	NA

Skip lines

read_csv(f, skip = 1)

A	B	C
1	2	3
4	5	NA

No header

read_csv(f, col_names = FALSE)

A	B	C
1	2	3

Read in a subset

read_csv(f, n_max = 1)

x	y	z
A	B	C
1	2	3
4	5	NA

Provide header

read_csv(f, col_names = c("x", "y", "z"))

A	B	C
NA	2	3
4	5	NA

Missing Values

read_csv(f, na = c("1", ""))

Read Non-Tabular Data

Read a file into a single string

read_file(file, locale = default_locale())

Read each line into its own string

read_lines(file, skip = 0, n_max = -1L, na = character(), locale = default_locale(), progress = interactive())

Read Apache style log files

read_log(file, col_names = FALSE, col_types = NULL, skip = 0, n_max = -1, progress = interactive())

Read a file into a raw vector

read_file_raw(file)

Read each line into a raw vector

read_lines_raw(file, skip = 0, n_max = -1L, progress = interactive())

Data types

readr functions guess the types of each column and convert types when appropriate (but will NOT convert strings to factors automatically).

A message shows the type of each column in the result.

```
## Parsed with column specification:
## cols(
##   age = col_integer(),
##   sex = col_character(),
##   earn = col_double()
## )
```

age is an integer

earn is a double (numeric)

sex is a character

1. Use **problems()** to diagnose problems.

x <- read_csv("file.csv"); problems(x)

2. Use a **col_** function to guide parsing.

- **col_guess()** - the default
- **col_character()**
- **col_double()**, **col_euro_double()**
- **col_datetime**(format = "") Also **col_date**(format = ""), **col_time**(format = "")
- **col_factor**(levels, ordered = FALSE)
- **col_integer()**
- **col_logical()**
- **col_number()**, **col_numeric()**
- **col_skip()**

x <- read_csv("file.csv", col_types = cols(
A = col_double(),
B = col_logical(),
C = col_factor()))

3. Else, read in as character vectors then parse with a **parse_** function.

- **parse_guess()**
- **parse_character()**
- **parse_datetime()** Also **parse_date()** and **parse_time()**
- **parse_double()**
- **parse_factor()**
- **parse_integer()**
- **parse_logical()**
- **parse_number()**

x\$A <- parse_number(x\$A)

Tidyverse cheatsheets!

Data Transformation with dplyr : : CHEAT SHEET

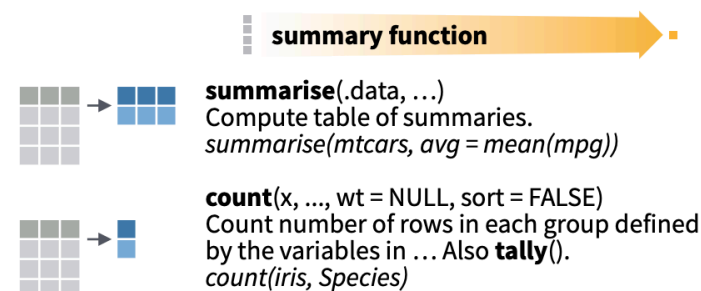


dplyr functions work with pipes and expect **tidy data**. In tidy data:



Summarise Cases

These apply **summary functions** to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).

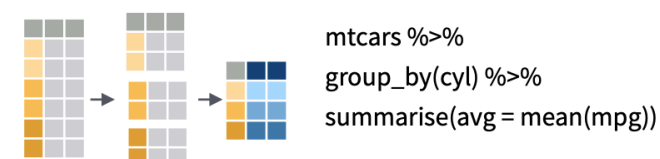


VARIATIONS

summarise_all() - Apply funs to every column.
summarise_at() - Apply funs to specific columns.
summarise_if() - Apply funs to all cols of one type.

Group Cases

Use **group_by()** to create a "grouped" copy of a table. dplyr functions will manipulate each "group" separately and then combine the results.



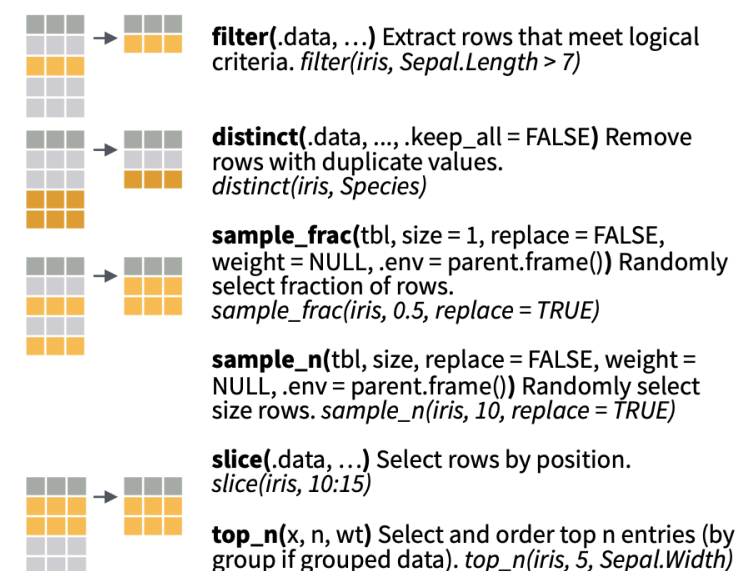
group_by(.data, ..., add = FALSE)
Returns copy of table grouped by ...
g_iris <- group_by(iris, Species)

ungroup(x, ...)
Returns ungrouped copy of table.
ungroup(g_iris)

Manipulate Cases

EXTRACT CASES

Row functions return a subset of rows as a new table.

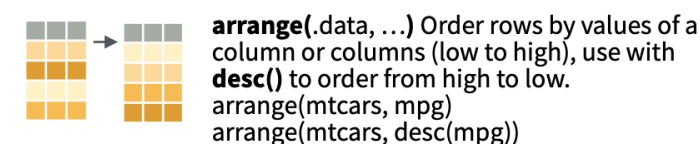


Logical and boolean operators to use with filter()

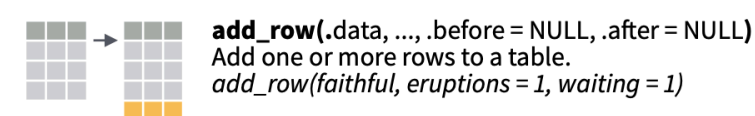
<	<=	is.na()	%in%		xor()
>	>=	!is.na()	!	&	

See **?base::Logic** and **?Comparison** for help.

ARRANGE CASES



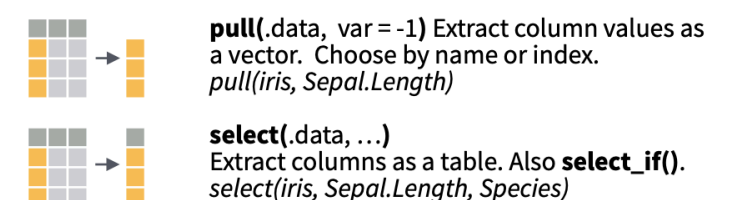
ADD CASES



Manipulate Variables

EXTRACT VARIABLES

Column functions return a set of columns as a new vector or table.

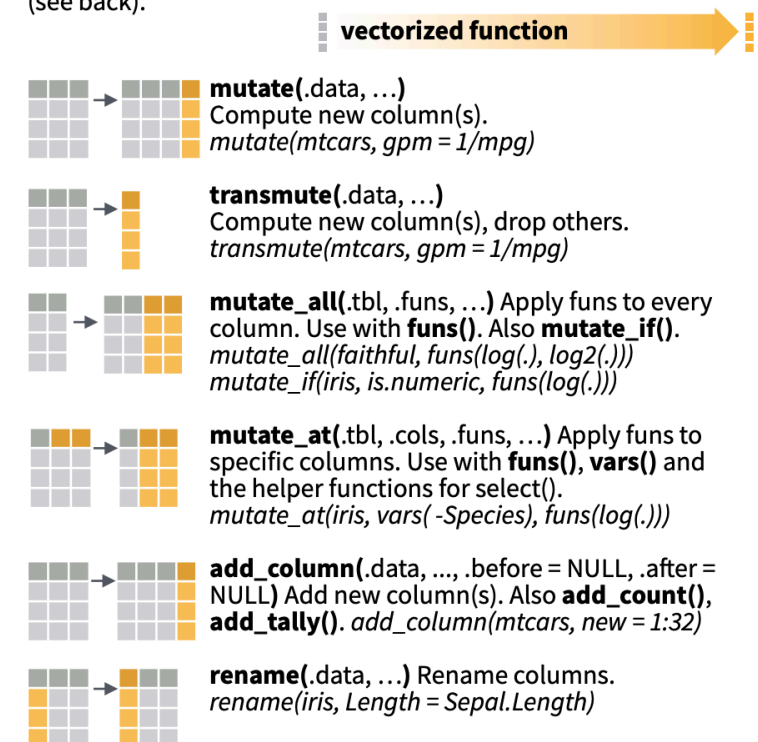


Use these helpers with **select()**, e.g. *select(iris, starts_with("Sepal"))*

contains(match)	num_range(prefix, range)	⋮, e.g. mpg:cyl
ends_with(match)	one_of(...)	⋮, e.g. -Species
matches(match)	starts_with(match)	

MAKE NEW VARIABLES

These apply **vectorized functions** to columns. Vectorized funs take vectors as input and return vectors of the same length as output (see back).



BREAK

After the break: Tutorial

Homework 1

Homework 1 - Data import (no Github yet)

- <https://github.com/tspringstein/259-langbasics-importing-hw>
- Fork this project
- Clone it to your RStudio
- Complete the homework (at least 4 completed questions counts as completed)
- Commit changes
- Push changes
- Invite Madison to be a collaborator on your project so she can provide comments

