



Research Data Workshop Series

- Introduction to R

Lucas Alcantara, Ph.D.
alcantal@uoguelph.ca



**AGRI-FOOD DATA
CANADA**

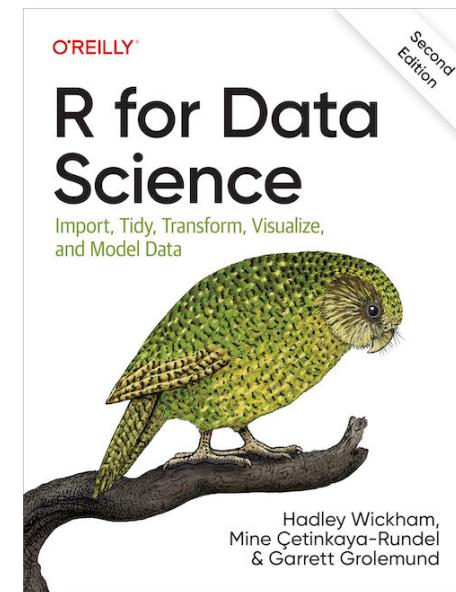
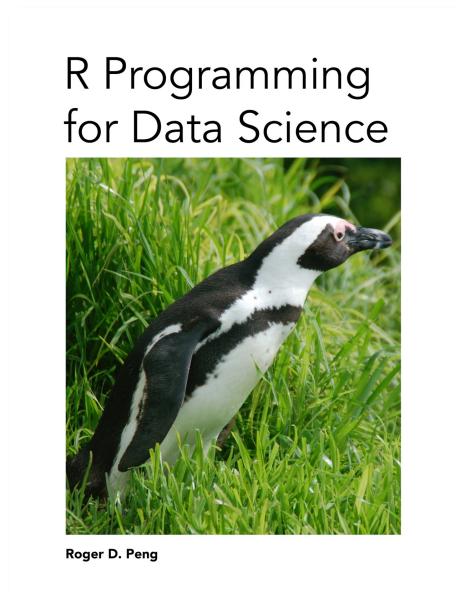
AT THE UNIVERSITY *of* GUELPH

Workshop Outline

- What is R
- Getting Started with R
- R Nuts and Bolts
- Getting Data in and Out of R
- Subsetting R Objects
- Project Time!
- Data Transformation with the Tidyverse

Based heavily on free e-books

- Peng, R. (2016). *R Programming for Data Science*. Lulu.com.
E-book available at: <https://bookdown.org/rdpeng/rprogdatascience>
- Wickham, H., Mine Cetinkaya-Rundel, & Grolemund, G. (2023). *R for data science: Import, tidy, transform, visualize, and model data* (2nd ed.). O'Reilly Media.
E-book available at: <https://r4ds.hadley.nz/>



What is R?

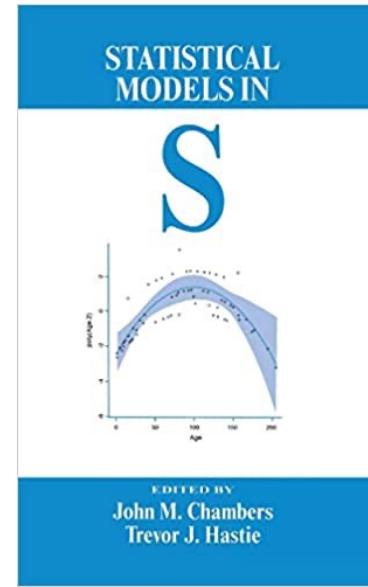
What is R?

The quick answer is: R is a dialect of S

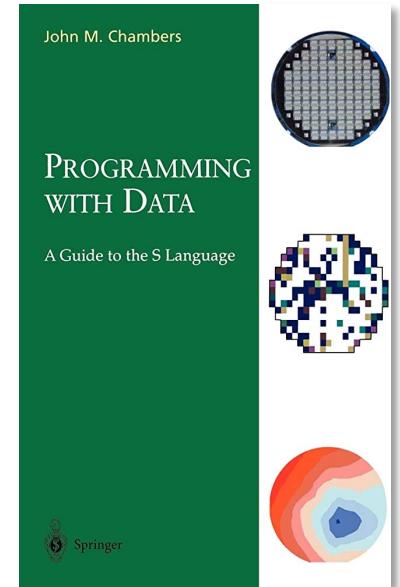
S is a language that was developed by John Chambers and others at the old Bell Telephone Laboratories (AT&T Corp.)

It was initiated in 1976 as an internal statistical analysis environment - originally implemented as *Fortran* libraries

In 1988, it was rewritten in *C* and began to resemble the system that we have today



1988



1998

The S Philosophy

“[W]e wanted *users* to be able to begin in an interactive environment, where they did not consciously think of themselves as programming. Then as their needs became clearer and their sophistication increased, they should be able to slide gradually into programming, when the language and system aspects would become more important.”

Stages in the Evolution of S, John Chambers

Back to R

- 1991: R was created by Ross Ihaka and Robert Gentleman in the Department of Statistics at the University of Auckland.
- 1993: The first announcement of R was made to the public.
- 1996: Paper published with author's experience developing R

Ross Ihaka and Robert Gentleman. *R: A language for data analysis and graphics*. Journal of Computational and Graphical Statistics, 5(3):299–314, 1996.
<https://doi.org/10.2307/1390807>

- 1997: R Core Group was formed, and it still controls the source code
- 2000: R version 1.0.0 was released
- 2023: R version 4.3.0 is the latest

R: A Language for Data Analysis and Graphics

ROSS IHAKA and Robert GENTLEMAN

In this article we discuss our experience designing and implementing a statistical computing language. In developing this new language, we sought to combine what we felt were useful features from two existing computer languages. We feel that the new language provides advantages in the areas of portability, computational efficiency, memory management, and scoping.

Key Words: Computer language; Statistical computing.

1. INTRODUCTION

This article discusses some issues involved in the design and implementation of a computer language for statistical data analysis. Our experience with these issues occurred while developing such a language. The work has been heavily influenced by two existing languages—Becker, Chambers, and Wilks' S (1985) and Steel and Sussman's Scheme (1975). We felt that there were strong points in each of these languages and that it would be interesting to see if the strengths could be combined. The resulting language is very similar in appearance to S, but the underlying implementation and semantics are derived from Scheme. In fact, we implemented the language by first writing an interpreter for a Scheme subset and then progressively mutating it to resemble S.

We added S-like features in several stages. First, we altered the language parser so that the syntax would resemble that of S. This created a major change in the appearance of the language, but it should be emphasized that the change was entirely superficial; the underlying semantics remained those of Scheme. Next, we modified the data types of the language by removing the single scalar data type we had put into our Scheme and replacing it with the vector-based types of S. This was a much more substantive change and required major modifications to the interpreter. The final substantive change involved adding the S notion of *lazy arguments* for functions.

At this point we had enough of a framework in place to begin building a full statistical language. This process is ongoing, but we feel that we are well on the way to building a complete and useful piece of software. The development of the key portions of language

Ross Ihaka is Senior Lecturer, and Robert Gentleman is Senior Lecturer, Department of Statistics, University of Auckland, Private Bag 92019, Auckland, New Zealand; e-mail: ihaka@stat.auckland.ac.nz.
©1996 American Statistical Association, Institute of Mathematical Statistics,
and Interface Foundation of North America
Journal of Computational and Graphical Statistics, Volume 5, Number 3, Pages 299-314

299



Design of the R System

The R system is divided into 2 conceptual parts:

- The “base” R system is available from CRAN: Comprehensive R Archive Network (<https://cran.r-project.org>).
- Everything else: many packages that can be used to extend the functionality of R.

The screenshot shows the homepage of The Comprehensive R Archive Network (CRAN). The top navigation bar includes links for "Download and Install R", "Source Code for all Platforms", and "Questions About R". The main content area features a large "R" logo and links to "CRAN", "Mirrors", "What's new?", "Search", and "CRAN Team". Below this are sections for "About R", "Software", "Documentation", and "FAQs". A prominent section titled "Download and Install R" provides links for precompiled binary distributions for Linux, macOS, and Windows. Another section, "Source Code for all Platforms", describes how to obtain source code for compilation. The "Questions About R" section contains frequently asked questions and links for answers. At the bottom, there is information about R and CRAN, and a note about mirrors.

Getting Started with R

Getting Started with R

Install R

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

The Comprehensive R Archive Network

Download and Install R

Precompiled binary distributions of the base system and contributed packages, **Windows and Mac** users most likely want one of these versions of R:

- [Download R for Linux \(Debian, Fedora/Redhat, Ubuntu\)](#)
- [Download R for macOS](#)
- [Download R for Windows](#)

R is part of many Linux distributions, you should check with your Linux package management system in addition to the link above.

Getting Started with R

Install R

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

R for Windows

Subdirectories:

base	Binaries for base distribution. This is what you want to install R for the first time .
contrib	Binaries of contributed CRAN packages (for R >= 3.4.x).
old contrib	Binaries of contributed CRAN packages for outdated versions of R (for R < 3.4.x).
Rtools	Tools to build R and R packages. This is what you want to build your own packages on Windows, or to build R itself.

Please do not submit binaries to CRAN. Package developers might want to contact Uwe Ligges directly in case of questions / suggestions related to Windows binaries.

You may also want to read the [R FAQ](#) and [R for Windows FAQ](#).

Note: CRAN does some checks on these binaries for viruses, but cannot give guarantees. Use the normal precautions with downloaded executables.



Getting Started with R

Install R

R-4.3.0 for Windows

[Download R-4.3.0 for Windows](#) (79 megabytes, 64 bit)

[README on the Windows binary distribution](#)

[New features in this version](#)

This build requires UCRT, which is part of Windows since Windows 10 and Windows Server 2016. On older systems, UCRT has to be installed manually from [here](#).

If you want to double-check that the package you have downloaded matches the package distributed by CRAN, you can compare the [md5sum](#) of the .exe to the [fingerprint](#) on the master server.

Frequently asked questions

- [Does R run under my version of Windows?](#)
- [How do I update packages in my previous version of R?](#)

Please see the [R FAQ](#) for general information about R and the [R Windows FAQ](#) for Windows-specific information.

Other builds

- Patches to this release are incorporated in the [r-patched snapshot build](#).
- A build of the development version (which will eventually become the next major release of R) is available in the [r-devel snapshot build](#).
- [Previous releases](#)

Note to webmasters: A stable link which will redirect to the current Windows binary release is
[<CRAN MIRROR>/bin/windows/base/release.html](#).

Last change: 2023-04-21

Getting Started with R

Install R

R for macOS

This directory contains binaries for the base distribution and of R and packages to run on macOS. R and package binaries for R versions older than 4.0.0 are only available from the [CRAN archive](#) so users of such versions should adjust the CRAN mirror setting (<https://cran-archive.r-project.org>) accordingly.

Note: Although we take precautions when assembling binaries, please use the normal precautions with downloaded executables.

R 4.3.0 "Already Tomorrow" released on 2023/04/21

Please check the integrity of the downloaded package by checking the signature:
`pkgutil --check-signature R-4.3.0.pkg`
in the *Terminal* application. If Apple tools are not available you can check the SHA1 checksum of the downloaded image:
`openssl sha1 R-4.3.0.pkg`

Latest release:

For Apple silicon (M1/M2) Macs:
[**R-4.3.0-arm64.pkg**](#)
SHA1-hash: 8ee0276daa9841993f218ebd2a8a7aa86c00d470
(ca. 90MB, notarized and signed)

For older Intel Macs:
[**R-4.3.0-x86_64.pkg**](#)
SHA1-hash: d28e528c8e3ec761aa4b871a8d444a1bfbee9bd3
(ca. 92MB, notarized and signed)

R 4.3.0 binary for macOS 11 (**Big Sur**) and higher, signed and notarized packages.
Contains R 4.3.0 framework, R.app GUI 1.79, Tcl/Tk 8.6.12 X11 libraries and Texinfo 6.8. The latter two components are optional and can be omitted when choosing "custom install", they are only needed if you want to use the `tcltk` R package or build package documentation from sources.

macOS Ventura users: there is a known bug in Ventura preventing installations from some locations without a prompt. If the installation fails, move the downloaded file away from the *Downloads* folder (e.g., to your home or Desktop).

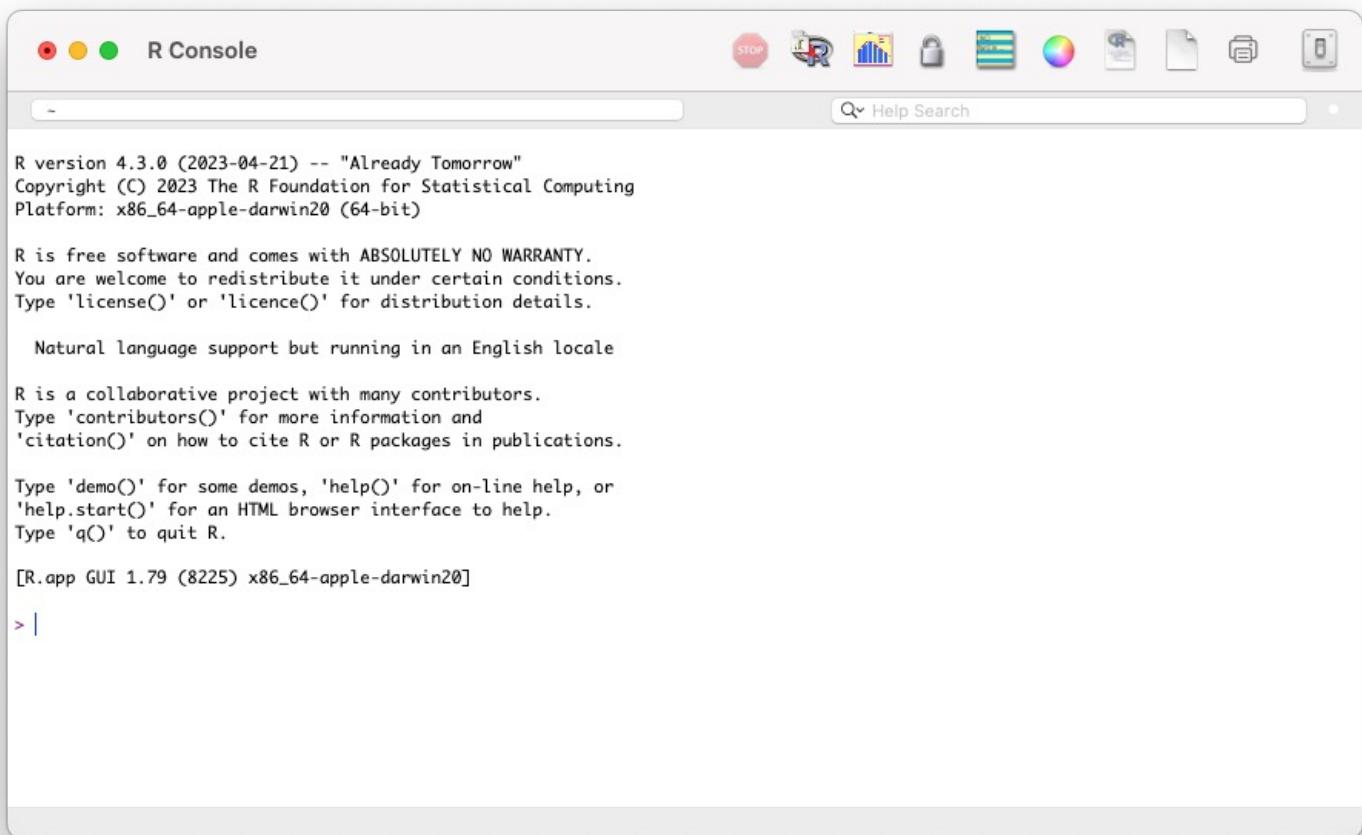
Note: the use of X11 (including `tcltk`) requires [XQuartz](#) (version 2.8.5 or later). Always re-install XQuartz when upgrading your macOS to a new major version.

This release uses Xcode 14.2/14.3 and GNU Fortran 12.2. If you wish to compile R packages which contain Fortran code, you may need to download the corresponding GNU Fortran compiler from <https://mac.R-project.org/tools>. Any external libraries and tools are expected to live in `/opt/R/arm64` (Apple silicon) or `/opt/R/x86_64` (Intel).

Or

Getting Started with R

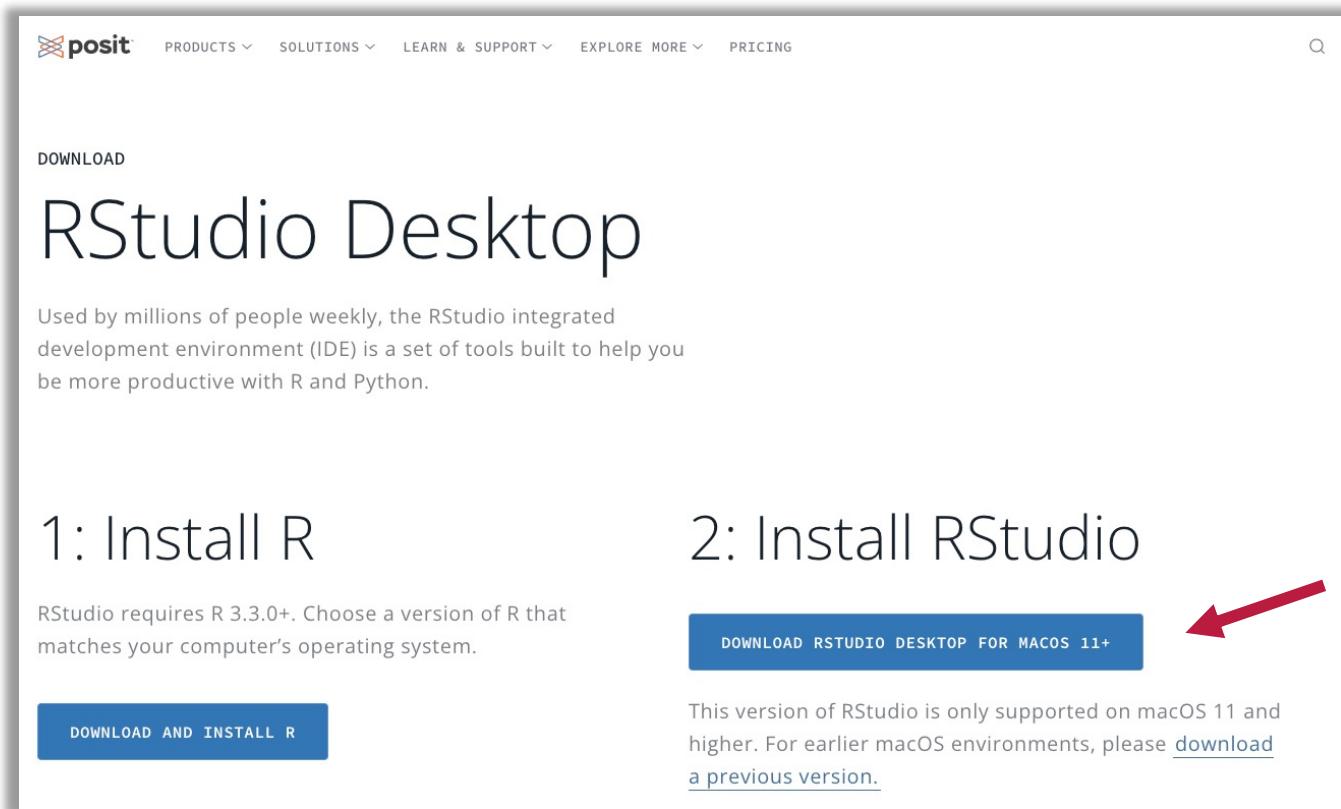
R Console: Not a very user-friendly interface



Getting Started with RStudio

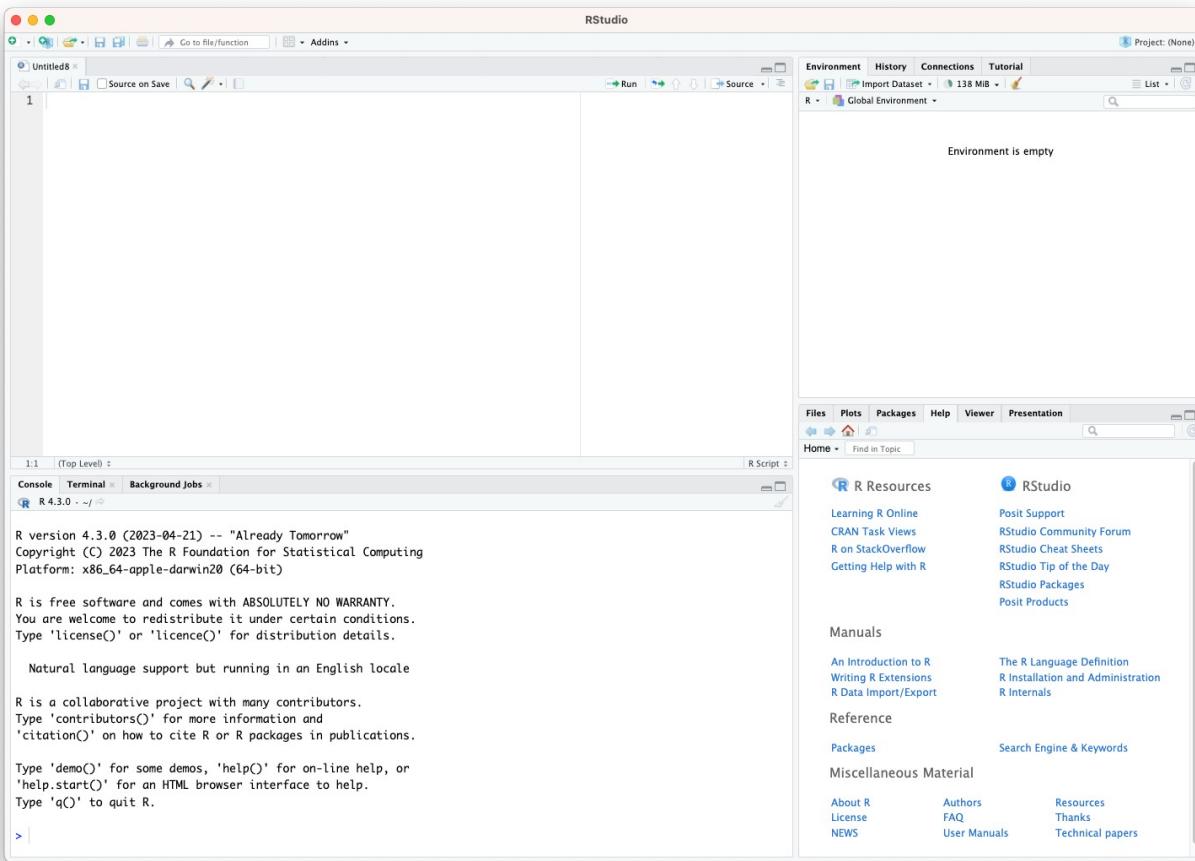
Install RStudio Desktop

<https://posit.co/download/rstudio-desktop>



Getting Started with RStudio

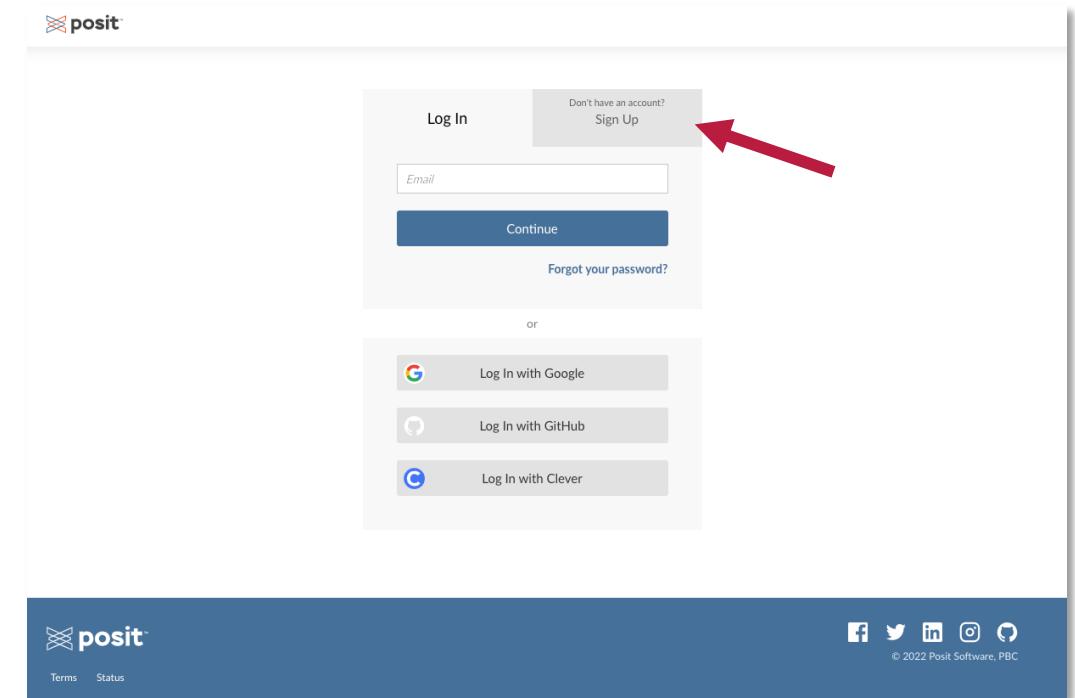
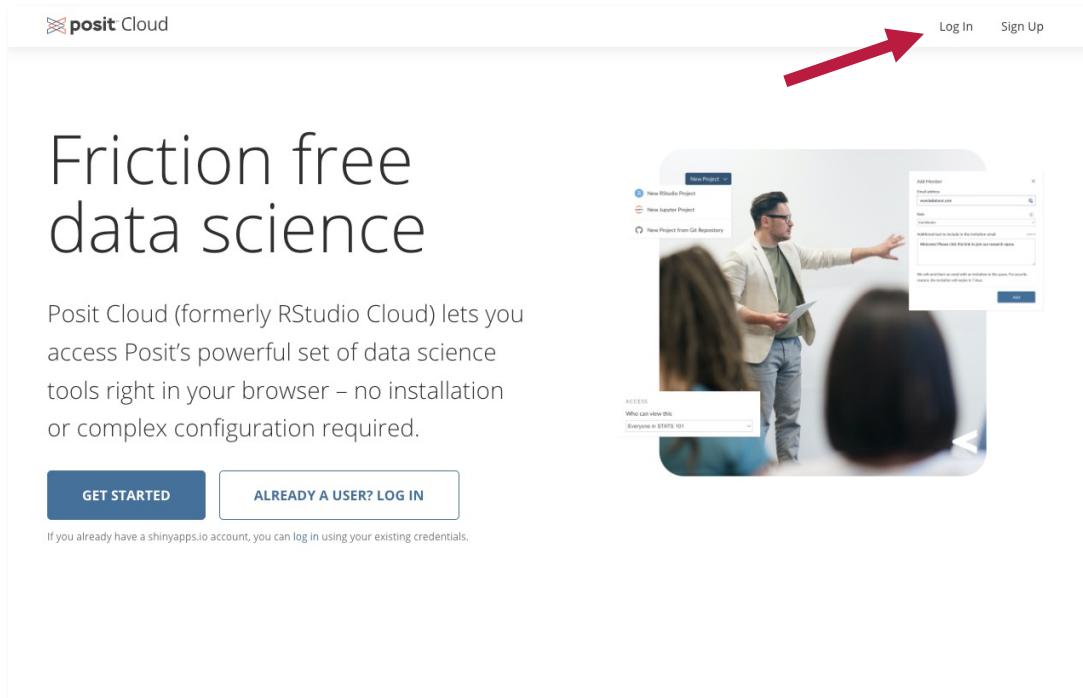
Install RStudio Desktop



Getting Started with RStudio

Another option: Posit Cloud (RStudio on the Cloud)

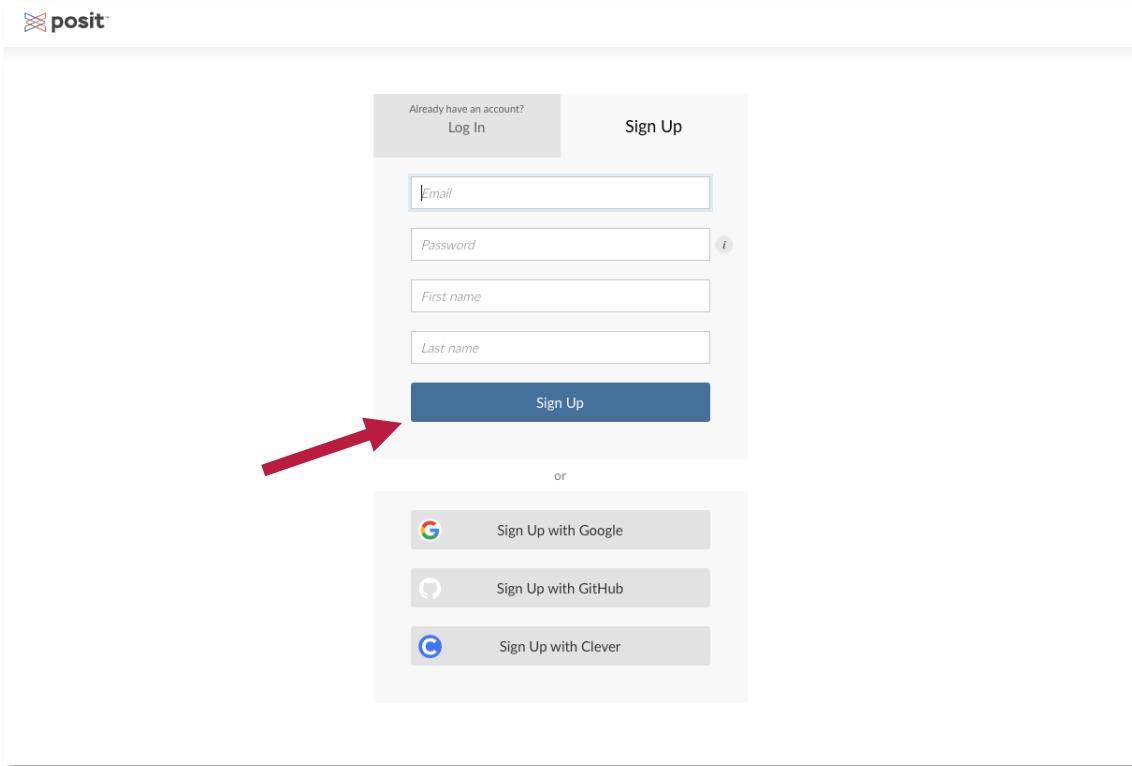
<https://posit.cloud>



Getting Started with RStudio

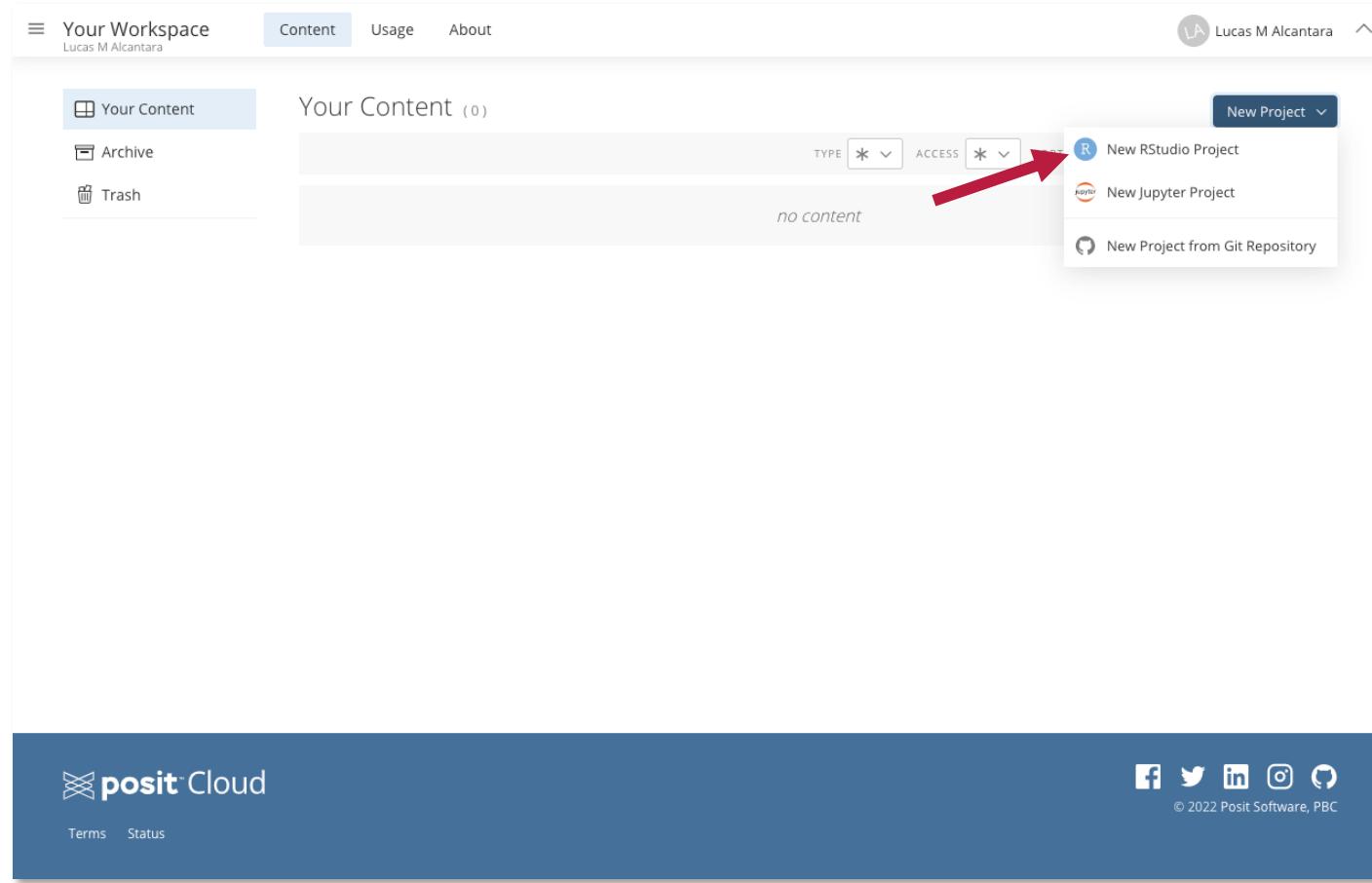
Another option: Posit Cloud (RStudio on the Cloud)

<https://posit.cloud>



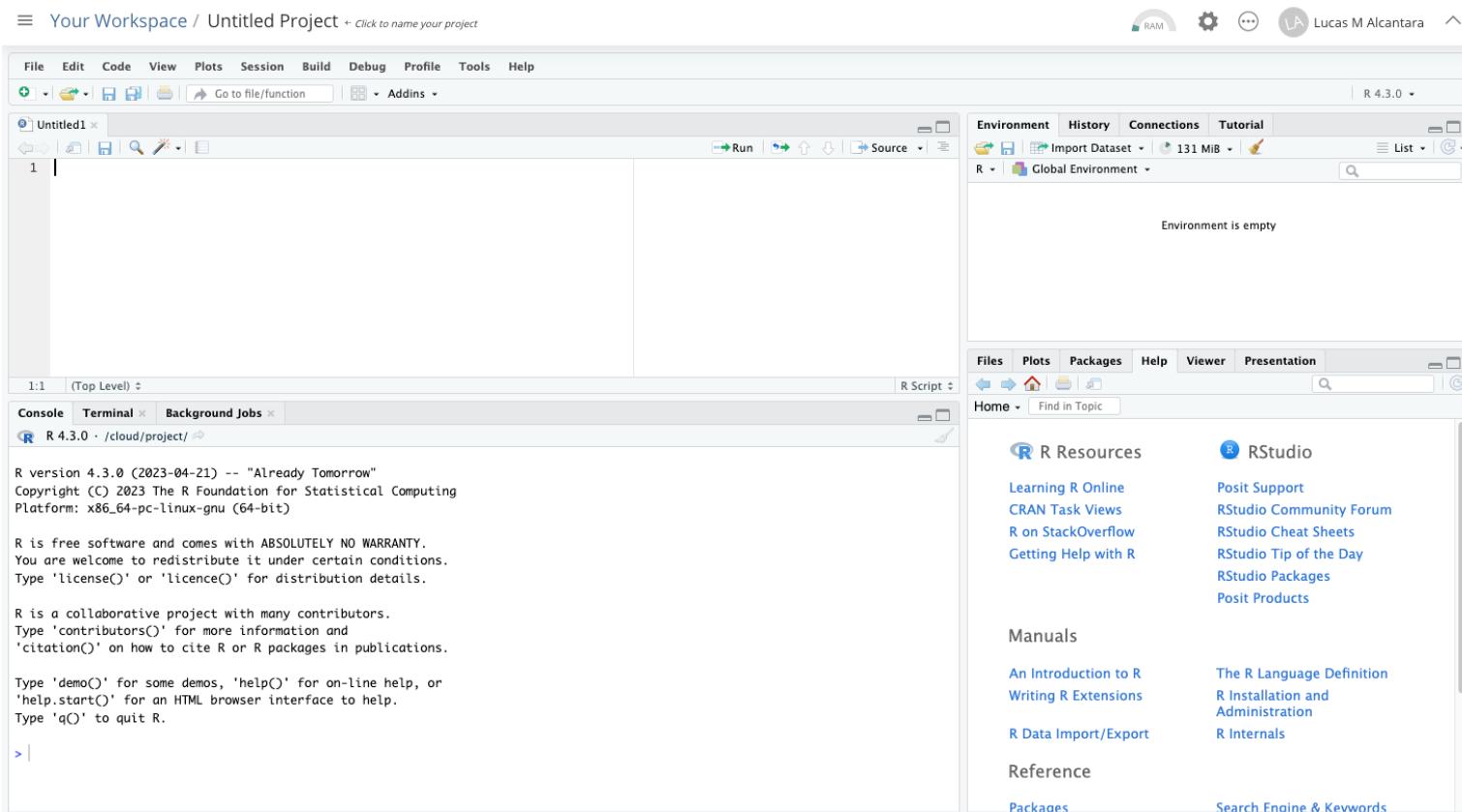
Getting Started with RStudio

Another option: Posit Cloud (RStudio on the Cloud)



Getting Started with RStudio

Another option: Posit Cloud (RStudio on the Cloud)



R Nuts and Bolts

Entering Input

We type expressions on the R console. The '<->' symbol is the assignment operator.

```
> x <- 1  
> print(x)  
[1] 1  
> x  
[1] 1  
> msg <- "hello"
```

The # character indicates a comment. Anything to the right of the # is ignored (including the # itself).

```
x <- ## Incomplete expression
```

Evaluation

When a complete expression is entered at the prompt, it is evaluated, and the result of the evaluated expression is returned. The result may be *auto-printed*.

```
> x <- 5 ## nothing printed  
> x      ## auto-printing occurs  
[1] 5  
> print(x) ## explicit printing  
[1] 5
```

The [1] shown on the output above indicates that 'x' is a vector and 5 is its first element.

```
> x <- 11:30  
> x  
[1] 11 12 13 14 15 16 17 18 19 20 21 22  
[13] 23 24 25 26 27 28 29 30
```

R Objects

R has five basic or “atomic” classes of objects:

- character
- numeric (real numbers)
- integer
- complex
- logical (True/False)

The most basic type of R object is a vector.

-> A vector can only contain objects of the same class.

Exception: list

Numbers

Numbers in R are generally treated as numeric objects, i.e., double precision real numbers

- $1 = 1.00$
- $2 = 2.00$

If you explicitly want an integer, you need to specify the 'L' suffix

- $1 = \text{numeric object}$
- $1L = \text{integer object}$

Numbers

There is also a special number 'Inf' which represents infinity

- $1/0 = \text{Inf}$
- $-1/0 = -\text{Inf}$
- $1/\text{Inf} = 0$

The value 'NaN' represents an undefined value ("not a number")

- $0 / 0 = \text{NaN}$

Attributes

Attributes are like metadata for the R object. If any, it can be accessed using the 'attributes()' function

- dimensions (e.g., matrices, arrays)
- class (e.g., integer, numeric)
- length
- other user-defined attributes/metadata
- etc.

Not all R objects contain attributes, in which case the attributes() function returns NULL.

Creating Vectors

The 'c()' function can be used to create vectors of objects by concatenating things together.

```
> x <- c(0.5, 0.6)      ## numeric  
> x <- c(TRUE, FALSE)    ## logical  
> x <- c(T, F)          ## logical  
> x <- c("a", "b", "c")  ## character  
> x <- 9:29               ## integer  
> x <- c(1+0i, 2+4i)     ## complex
```

Mixing Objects

What will be the class of y on each of the following codes?

```
> y <- c(1.7, "a")
> y <- c(TRUE, 2)
> y <- c("a", TRUE)
```

Mixing Objects

What will be the class of y on each of the following codes?

```
> y <- c(1.7, "a")    ## character  
> y <- c(TRUE, 2)     ## numeric  
> y <- c("a", TRUE)   ## character
```

When different objects are mixed in a vector, coercion occurs so that every element in the vector is of the same class.

Explicit Coercion

Objects can be explicitly coerced from one class to another using the 'as.*' functions, if available.

```
> x <- 0:6
> class(x)
[1] "integer"
> as.numeric(x)
[1] 0 1 2 3 4 5 6
> as.logical(x)
[1] FALSE TRUE TRUE TRUE TRUE TRUE TRUE
> as.character(x)
[1] "0" "1" "2" "3" "4" "5" "6"
```

Explicit Coercion

When R can't figure out how to coerce an object, it can result in NAs being introduced by coercion. A warning message is usually shown by R:

```
> x <- c("a", "b", "c")
> as.numeric(x)
Warning: NAs introduced by coercion
[1] NA NA NA
> as.logical(x)
[1] NA NA NA
> as.complex(x)
Warning: NAs introduced by coercion
[1] NA NA NA
```

Matrices

Matrices are vectors with a dimension attribute. The dimension attribute is itself an integer vector of length 2 (number of rows, number of columns)

```
> m <- matrix(nrow = 2, ncol = 3)
> m
 [,1] [,2] [,3]
[1,] NA   NA   NA
[2,] NA   NA   NA
> dim(m)
[1] 2 3
> attributes(m)
$dim
[1] 2 3
```

Matrices

Matrices are constructed *column-wise*, so entries can be thought of starting in the “upper left” corner (1,1) and running down the columns (1,2; 2,1; 2,2...).

```
> m <- matrix(1:6, nrow = 2, ncol = 3)
> m
 [,1] [,2] [,3]
[1,]    1    3    5
[2,]    2    4    6
```

Matrices

Matrices can also be created directly from vectors by adding a dimension attribute.

```
> m <- 1:10  
> m  
[1] 1 2 3 4 5 6 7 8 9 10  
> dim(m) <- c(2, 5)  
> m  
[,1] [,2] [,3] [,4] [,5]  
[1,] 1 3 5 7 9  
[2,] 2 4 6 8 10
```

Matrices

Matrices can be created by column-binding or row-binding with the 'cbind()' and 'rbind()' functions.

```
> x <- 1:3  
> y <- 10:12  
> cbind(x, y)  
      x   y  
[1,] 1 10  
[2,] 2 11  
[3,] 3 12  
> rbind(x, y)  
 [,1] [,2] [,3]  
x     1     2     3  
y    10    11    12
```

Lists

Lists are a special type of vector that can contain elements of different classes.

```
> x <- list(1, "a", TRUE, 1 + 4i)
> x
[[1]]
[1] 1

[[2]]
[1] "a"

[[3]]
[1] TRUE

[[4]]
[1] 1+4i
```

Lists

Lists are a special type of vector that can contain elements of different classes.

```
> x <- list(1, "a", TRUE, 1 + 4i)  
  
> x  
[[1]]  
[1] 1  
  
[[2]]  
[1] "a"  
  
[[3]]  
[1] TRUE  
  
[[4]]  
[1] 1+4i
```

Based on the functions you learned, how can you tell what the classes of each of these objects are?

Factors

Factors are used to represent categorical data and can be unordered or ordered. “Integer vector labelled in alphabetical order”

```
> x <- factor(c("yes", "yes", "no", "yes", "no"))
> x
[1] yes yes no  yes no
Levels: no yes
> table(x)
x
no yes
 2   3
> ## See the underlying representation of factor
> unclass(x)
[1] 2 2 1 2 1
attr(,"levels")
[1] "no"  "yes"
```

Factors

The order of the levels of a factor can be set using the ‘levels’ argument to `factor()`.

```
> x <- factor(c("yes", "yes", "no", "yes", "no"))
> x ## Levels are put in alphabetical order
[1] yes yes no  yes no
Levels: no yes
> x <- factor(c("yes", "yes", "no", "yes", "no"),
+               levels = c("yes", "no"))
> x
[1] yes yes no  yes no
Levels: yes no
```

Missing Values

Missing values are denoted by NA or NaN for q undefined mathematical operations.

- `is.na()` is used to test if objects are NA
- `is.nan()` is used to test if objects are NaN
- NA values have a class (integer NA, character NA, etc.)
- NaN values are NA, but NOT the other way around

Missing Values

```
> ## Create a vector with NAs in it
> x <- c(1, 2, NA, 10, 3)
> ## Return a logical vector indicating which elements are NA
> is.na(x)
[1] FALSE FALSE TRUE FALSE FALSE
> ## Return a logical vector indicating which elements are NaN
> is.nan(x)
[1] FALSE FALSE FALSE FALSE FALSE

> ## Now create a vector with both NA and NaN values
> x <- c(1, 2, NaN, NA, 4)
> is.na(x)
[1] FALSE FALSE TRUE TRUE FALSE
> is.nan(x)
[1] FALSE FALSE TRUE FALSE FALSE
```

Data Frames

Data frames are used to store tabular data in R. It's a special type of list where every element of the list must have the same length.

- Columns = Elements of the list
- Number of rows = Length of each element of the list

Data frames can store different classes of objects in each column.

Data frames have both column and row names

Data Frames

Data frames can be created using the 'data.frame()' function

```
> x <- data.frame(foo = 1:4, bar = c(T, T, F, F))

> x
  foo   bar
1 1  TRUE
2 2  TRUE
3 3 FALSE
4 4 FALSE

> nrow(x)
[1] 4

> ncol(x)
[1] 2
```

Names

R objects can have names, which is very useful for writing readable code and self-describing objects.

```
> x <- 1:3
> names(x)
NULL
> names(x) <- c("New York", "Seattle", "Los Angeles")
> x
  New York      Seattle Los Angeles
    1           2           3
> names(x)
[1] "New York"     "Seattle"      "Los Angeles"
```

Names

R objects can have names, which is very useful for writing readable code and self-describing objects.

```
> x <- list("Los Angeles" = 1, Boston = 2, London = 3)
> x
$`Los Angeles`
[1] 1

$Boston
[1] 2

$London
[1] 3
> names(x)
[1] "Los Angeles" "Boston"      "London"
```

Names

R objects can have names, which is very useful for writing readable code and self-describing objects.

```
> m <- matrix(1:4, nrow = 2, ncol = 2)
> dimnames(m) <- list(c("a", "b"), c("c", "d"))
> m
   c d
a 1 3
b 2 4

> colnames(m) <- c("h", "f")
> rownames(m) <- c("x", "z")
> m
   h f
x 1 3
z 2 4
```

Getting Data in and Out of R

Reading and Writing Data

There are a few ways to read and write data into R.

Target	Read function	Write Function
Tabular data	<code>read.table()</code> , <code>read.csv()</code> , etc.	<code>write.table()</code> , <code>write.csv()</code> , etc.
Lines	<code>readLines()</code>	<code>writeLines()</code>
R code files	<code>source()</code>	-
Workspace files	<code>load()</code>	<code>save()</code>

Reading Data Files with `read.table()`

The `read.table()` function is one of the most used functions for reading data.

Let's look at the Help page: `?read.table`

Reading Data Files with `read.table()`

A few tips on what you can do to read your data faster:

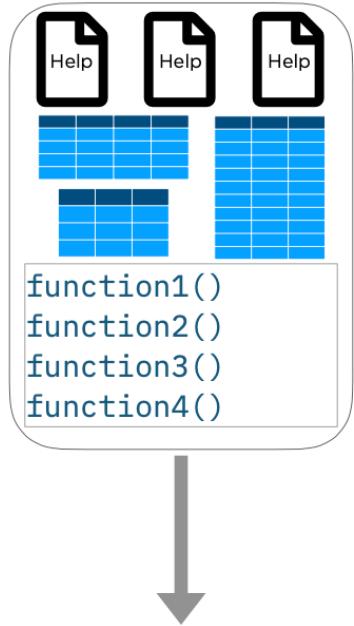
- Set the argument `comment.char = ""`
- Set the argument `stringsAsFactors = FALSE`
- Use the argument `colClasses`
 - Give the specific classes

```
> initial <- read.table("datatable.txt", nrows = 100)
> classes <- sapply(initial, class)
> tabAll <- read.table("datatable.txt", colClasses = classes)
```

- Assume all columns have the same class, e.g., character
- Use faster functions from different packages

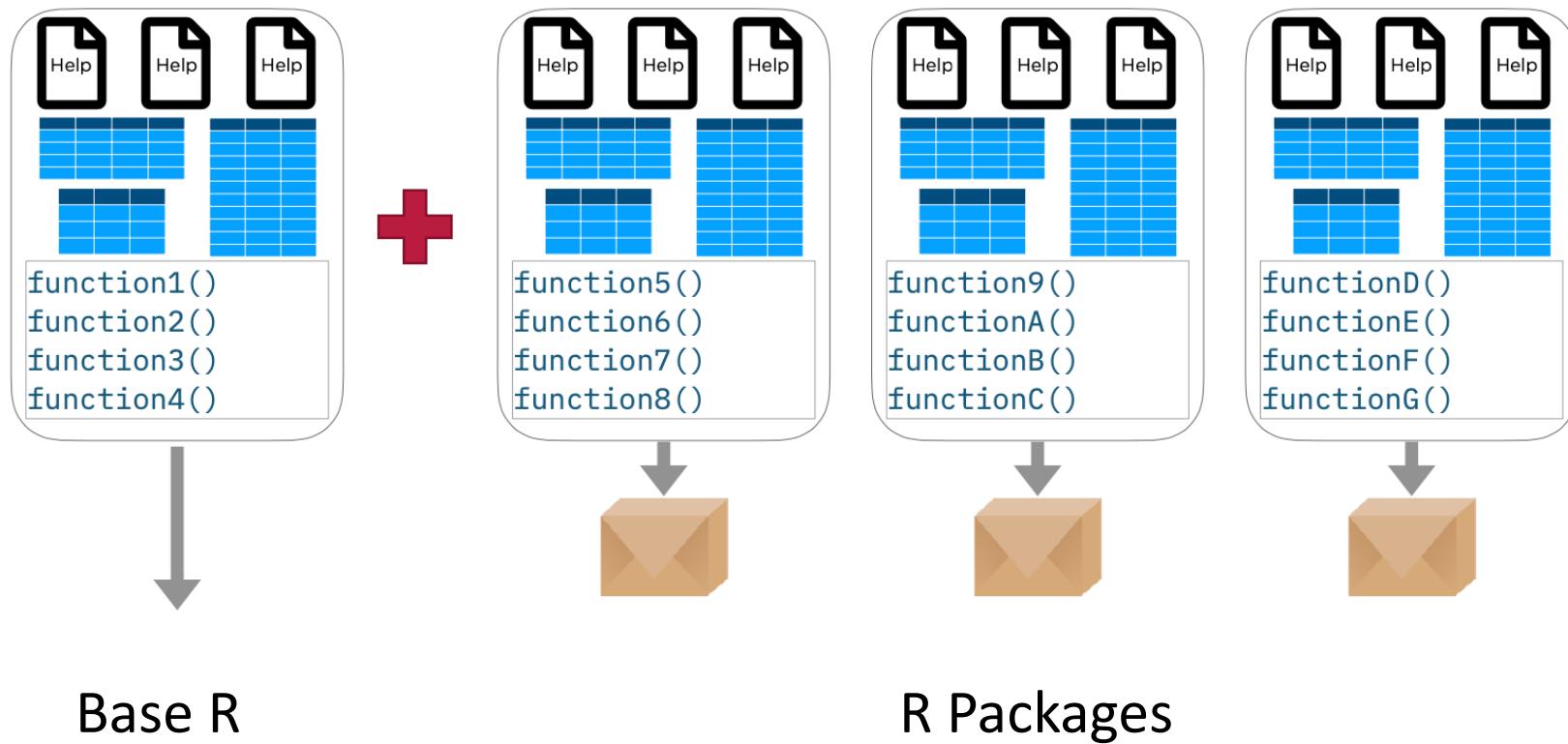
R Packages

R Packages



Base R

R Packages



R Packages on the Console

Install a package: once per computer/user

- `install.packages("readr")`

Load a package: once per R session

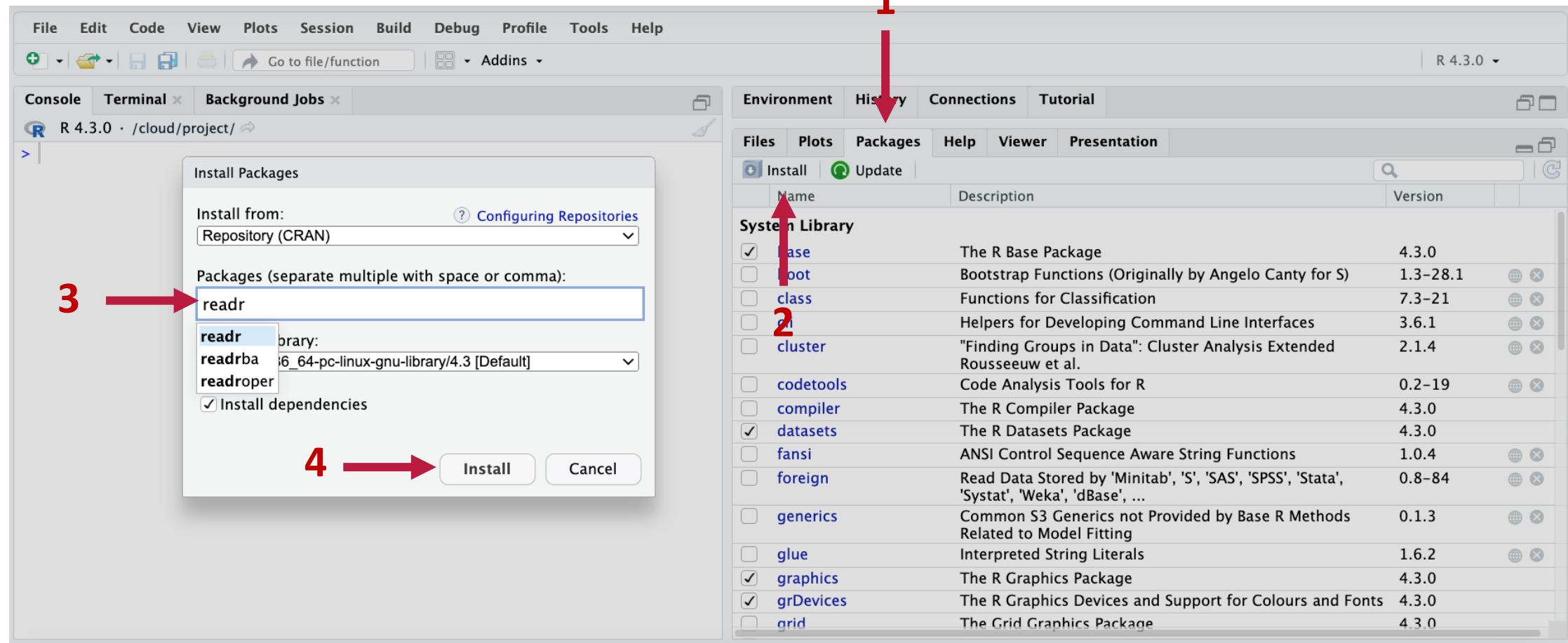
- `library(readr)`

Remove a package: once per computer/user

- `remove.packages("readr")`

R Packages on RStudio

Install a package: once per computer/user



R Packages on RStudio

Load a package: once per R session

The screenshot shows the RStudio interface with the following details:

- Console Output:** Shows the results of running `install.packages` to install several packages. The output includes:
 - * installing *binary* package 'clipr' ...
 - * DONE (clipr)
 - * installing *binary* package 'crayon' ...
 - * DONE (crayon)
 - * installing *binary* package 'hms' ...
 - * DONE (hms)
 - * installing *binary* package 'cpp11' ...
1.
 - * DONE (cpp11)
 - * installing *binary* package 'bit64' ...
 - * DONE (bit64)
 - * installing *binary* package 'progress' ...
 - * DONE (progress)
 - * installing *binary* package 'tzdb' ...
 - * DONE (tzdb)
 - * installing *binary* package 'vroom' ...
 - * DONE (vroom)
 - * installing *binary* package 'readr' ...
 - * DONE (readr)

The message "The downloaded source packages are in '/tmp/RtmpYEBH3j/downloaded_packages'" is also present.
- Packages Tab:** Shows the installed packages. The 'stats' package is selected (indicated by a checked checkbox). Other packages listed include progress, R6, readr, rlang, rpart, spatial, splines, survival, tcltk, tibble, tidyselect, and tools.

R Packages on RStudio

Load a package: once per R session

The screenshot shows the RStudio interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, and Help. The version R 4.3.0 is displayed. The left sidebar has tabs for Console, Terminal (active), Background Jobs, and a file icon. The Terminal tab shows the command line output:

```
R 4.3.0 · /cloud/project/
* installing *binary* package 'clipr' ...
* DONE (clipr)
* installing *binary* package 'crayon' ...
* DONE (crayon)
* installing *binary* package 'hms' ...
* DONE (hms)
* installing *binary* package 'cpp11' ...
* DONE (cpp11)
* installing *binary* package 'bit64' ...
* DONE (bit64)
* installing *binary* package 'progress' ...
* DONE (progress)
* installing *binary* package 'tzdb' ...
* DONE (tzdb)
* installing *binary* package 'vroom' ...
* DONE (vroom)
* installing *binary* package 'readr' ...
* DONE (readr)

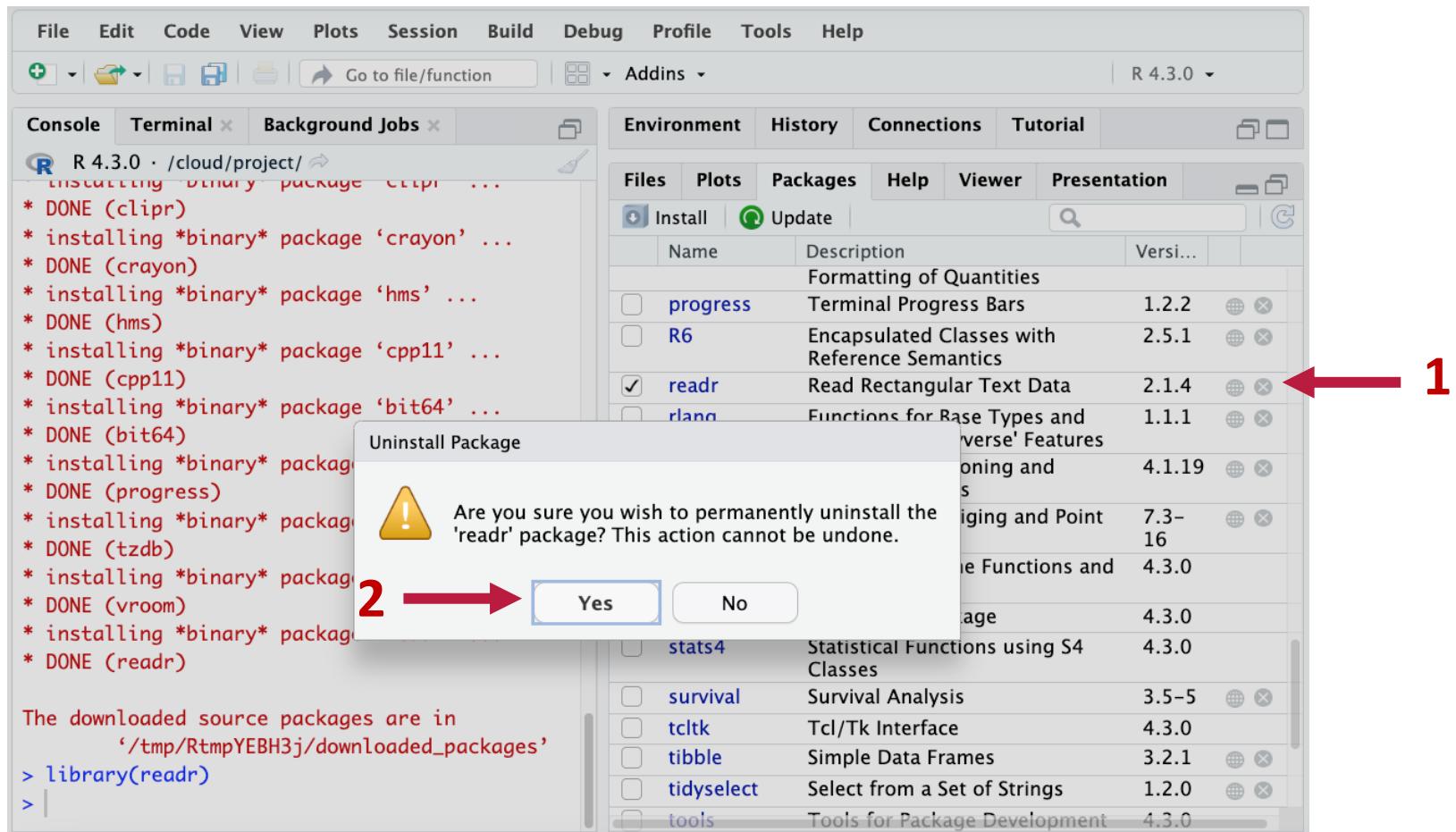
The downloaded source packages are in
  '/tmp/RtmpYEBH3j/downloaded_packages'
> library(readr)
>
```

The right panel features tabs for Environment, History, Connections, Tutorial, Files, Plots, Packages (selected), Help, Viewer, and Presentation. The Packages tab displays a list of installed packages:

Name	Description	Vers...
progress	Terminal Progress Bars	1.2.2
R6	Encapsulated Classes with Reference Semantics	2.5.1
<input checked="" type="checkbox"/> readr	Read Rectangular Text Data	2.1.4
<input type="checkbox"/> rlang	Functions for Base Types and Core R and 'Tidyverse' Features	1.1.1
<input type="checkbox"/> rpart	Recursive Partitioning and Regression Trees	4.1.19
<input type="checkbox"/> spatial	Functions for Kriging and Point Pattern Analysis	7.3-16
<input type="checkbox"/> splines	Regression Spline Functions and Classes	4.3.0
<input checked="" type="checkbox"/> stats	The R Stats Package	4.3.0
<input type="checkbox"/> stats4	Statistical Functions using S4 Classes	4.3.0
<input type="checkbox"/> survival	Survival Analysis	3.5-5
<input type="checkbox"/> tcltk	Tcl/Tk Interface	4.3.0
<input type="checkbox"/> tibble	Simple Data Frames	3.2.1
<input type="checkbox"/> tidyselect	Select from a Set of Strings	1.2.0
<input type="checkbox"/> tools	Tools for Package Development	4.3.0

R Packages on RStudio

Remove a package: once per computer/user



R Packages on RStudio

Remove a package: once per computer/user

The screenshot shows the RStudio interface with the following details:

- Console Tab:** Displays the command-line output of package installations and removals.
- Packages Tab:** Shows the list of installed packages in the System Library.

Console Output:

```
R 4.3.0 · /cloud/project/ ↵
  * INSTALLING "binary" package 'rims' ...
  * DONE (rims)
  * installing *binary* package 'cpp11' ...
  * DONE (cpp11)
  * installing *binary* package 'bit64' ...
  * DONE (bit64)
  * installing *binary* package 'progress' ...
  * DONE (progress)
  * installing *binary* package 'tzdb' ...
  * DONE (tzdb)
  * installing *binary* package 'vroom' ...
  * DONE (vroom)
  * installing *binary* package 'readr' ...
  * DONE (readr)

The downloaded source packages are in
  '/tmp/RtmpYEBH3j/downloaded_packages'
> library(readr)
> remove.packages("readr")
Removing package from '/cloud/lib/x86_64-pc-linux-gnu-library/4.3'
(as 'lib' is unspecified)
>
```

Packages Tab: Shows the System Library with the following packages listed:

Name	Description	Version
base	The R Base Package	4.3.0
bit	Classes and Methods for Fast Memory-Efficient Boolean Selections	4.0.5
bit64	A S3 Class for Vectors of 64bit Integers	4.0.5
boot	Bootstrap Functions (Originally by Angelo Canty for S)	1.3-28.1
class	Functions for Classification	7.3-21
cli	Helpers for Developing Command Line Interfaces	3.6.1
clipr	Read and Write from the System Clipboard	0.8.0
cluster	"Finding Groups in Data": Cluster Analysis Extended Rousseeuw et al.	2.1.4
codetools	Code Analysis Tools for R	0.2-19
compiler	The R Compiler Package	4.3.0
rnn11	A C++11 Interface for R's C	0.4.3

Conflicting R Packages

Different/Similar functions might have the same name in different packages

- Install and load the ‘dplyr’ package
- Same function name, different packages
- Load order DO matter!
 - Detach both ‘stats’ and ‘dplyr’ packages
 - Load ‘dplyr’ first and ‘stats’ later. Did you notice any difference?
- Try to be specific to which package you want to use to avoid conflicts
 - `stats::filter()`
 - `dplyr::filter()`

TIP: Use the ‘conflicted’ package to manage conflicts

R Packages

Loading a package silently:

- `suppressPackageStartupMessages(library(dplyr))`
- Use it with caution!

Getting Data in and Out of R

Reading Data Files with the `readr` package

Popular read functions: `read_table()`, `read_csv()`, `read_fwf`

Advantage of using `readr` functions to read data files

- Easier debugging: warnings indicate which rows/observations triggered them
- Faster reading: automatically guesses column types from the first 1k lines only
- Reads compressed files automatically
- Nice user-oriented features:
 - Progress bar when reading big files
 - Short data description
- R objects from `readr` are tibbles, not base R's data frames

*More on this later

Reading Multiple Data Files

Sometimes we have data saved in multiple files (e.g., one file per day)

What's a quick way to read all of them into one R Object?

```
1 # Load data.table package
2 library(data.table)
3
4 # Create a list of files in a directory
5 file_list <- list.files(".")
6
7 # Apply the fread function to each of those files and store in a list
8 data_list <- lapply(file_list, fread)
9
10 # Bind rows in a list
11 data <- data.table::rbindlist(data_list)
```

*Use data from the insentec folder

Writing Data Out of R

Like reading, many functions are available to write data out of R. A few examples:

- `write.csv | write_csv`
- `write.table | write_table`
- `fwrite`

Most of them you start specifying the object followed by the file path:

- `write.csv(data, "my/path/to/data.csv")`

Common arguments include:

- Column/row names (T/F); Quote (T/F); Encoding (utf8, latin-1, etc.);
- Field separator (comma, space, tab, pipe, etc.); String for NA values

Subsetting R Objects

Common Subsetting Operators

There are three operators that can be used to extract subsets of R objects.

- The `[` operator
 - Always returns an object of the same class as the original
 - It can be used to select multiple elements of an object
- The `[[` operator
 - Is used to extract elements of a list or a data frame
 - It can only be used to extract a single element
 - The class of the returned object will not necessarily be a list or data frame
- The `$` operator
 - Is used to extract elements of a list or data frame by literal name
 - Its semantics are similar to `[[`

Subsetting a Vector

The [operator can be used to extract multiple elements of a vector by passing the operator an integer, an integer sequence, or a logical sequence.

```
> x <- c("a", "b", "c", "c", "d", "a")
> x[1]
[1] "a"
> x[2]
[1] "b"
> x[1:4]
[1] "a" "b" "c" "c"
> x[c(1, 3, 4)]
[1] "a" "c" "c"
> x[x > "a"]
[1] "b" "c" "c" "d"
```

Subsetting a Matrix

Matrices can be subsetted in the usual way with (i,j) type indices. One of the indices can be missing.

```
> x <- matrix(1:6, 2, 3)
> x
     [,1] [,2] [,3]
[1,]    1    3    5
[2,]    2    4    6
> x[1, 2]
[1] 3
> x[2, 1]
[1] 2
> x[1, ] ## Extract the first row
[1] 1 3 5
> x[, 2] ## Extract the second column
[1] 3 4
> x[1, , drop = FALSE]
     [,1] [,2] [,3]
[1,]    1    3    5
```

Subsetting a List

Lists in R can be subsetted using all three operators, each used for a different purpose.

```
> x <- list(foo = 1:4, bar = 0.6)
> x
$foo
[1] 1 2 3 4

$bar
[1] 0.6

> x[[1]]
[1] 1 2 3 4
> x[["bar"]]
[1] 0.6
> x$bar
[1] 0.6
```

*Notice you don't need the quotes when you use the \$ operator.

Subsetting a List

The \$ operator can only be used with literal names.

```
> x <- list(foo = 1:4, bar = 0.6, baz = "hello")  
> name <- "foo"  
>  
> ## computed index for "foo"  
> x[[name]]  
[1] 1 2 3 4  
>  
> ## element "name" doesn't exist! (but no error here)  
> x$name  
NULL  
>  
> ## element "foo" does exist  
> x$foo  
[1] 1 2 3 4
```

Subsetting Nested Elements of a List

The `[]` operator can take an integer sequence if you want to extract a nested element.

```
> x <- list(a = list(10, 12, 14), b = c(3.14, 2.81))  
>  
> ## Get the 3rd element of the 1st element  
> x[[c(1, 3)]]  
[1] 14  
>  
> ## Same as above  
> x[[1]][[3]]  
[1] 14  
>  
> ## 1st element of the 2nd element  
> x[[c(2, 1)]]  
[1] 3.14
```

Subsetting Multiple Elements of a List

The [operator can be used to extract multiple elements from a list.

```
> x <- list(foo = 1:4, bar = 0.6, baz = "hello")
> x[c(1, 3)]
$foo
[1] 1 2 3 4

$baz
[1] "hello"
```

*Note that `x[c(1, 3)]` is NOT the same as `x[[c(1, 3)]]`.

Removing NA Values

A common task in data analysis is removing missing values (NAs).

```
> x <- c(1, 2, NA, 4, NA, 5)
> bad <- is.na(x)
> print(bad)
[1] FALSE FALSE  TRUE FALSE  TRUE FALSE
> x[!bad]
[1] 1 2 4 5
```

Removing NA Values

What if there are multiple R objects and you want to take the subset with no missing values in any of those objects?

```
> x <- c(1, 2, NA, 4, NA, 5)
> y <- c("a", "b", NA, "d", NA, "f")
> good <- complete.cases(x, y)
> good
[1] TRUE TRUE FALSE TRUE FALSE TRUE
> x[good]
[1] 1 2 4 5
> y[good]
[1] "a" "b" "d" "f"
```

*Note that both vectors must have the same lengths

Removing NA Values

We can also use `complete.cases()` with data frames

```
> head(airquality)
   Ozone Solar.R Wind Temp Month Day
1    41     190  7.4   67     5   1
2    36     118  8.0   72     5   2
3    12     149 12.6   74     5   3
4    18     313 11.5   62     5   4
5    NA      NA 14.3   56     5   5
6    28      NA 14.9   66     5   6
> good <- complete.cases(airquality)
> head(airquality[good, ])
   Ozone Solar.R Wind Temp Month Day
1    41     190  7.4   67     5   1
2    36     118  8.0   72     5   2
3    12     149 12.6   74     5   3
4    18     313 11.5   62     5   4
7    23     299  8.6   65     5   7
8    19      99 13.8   59     5   8
```

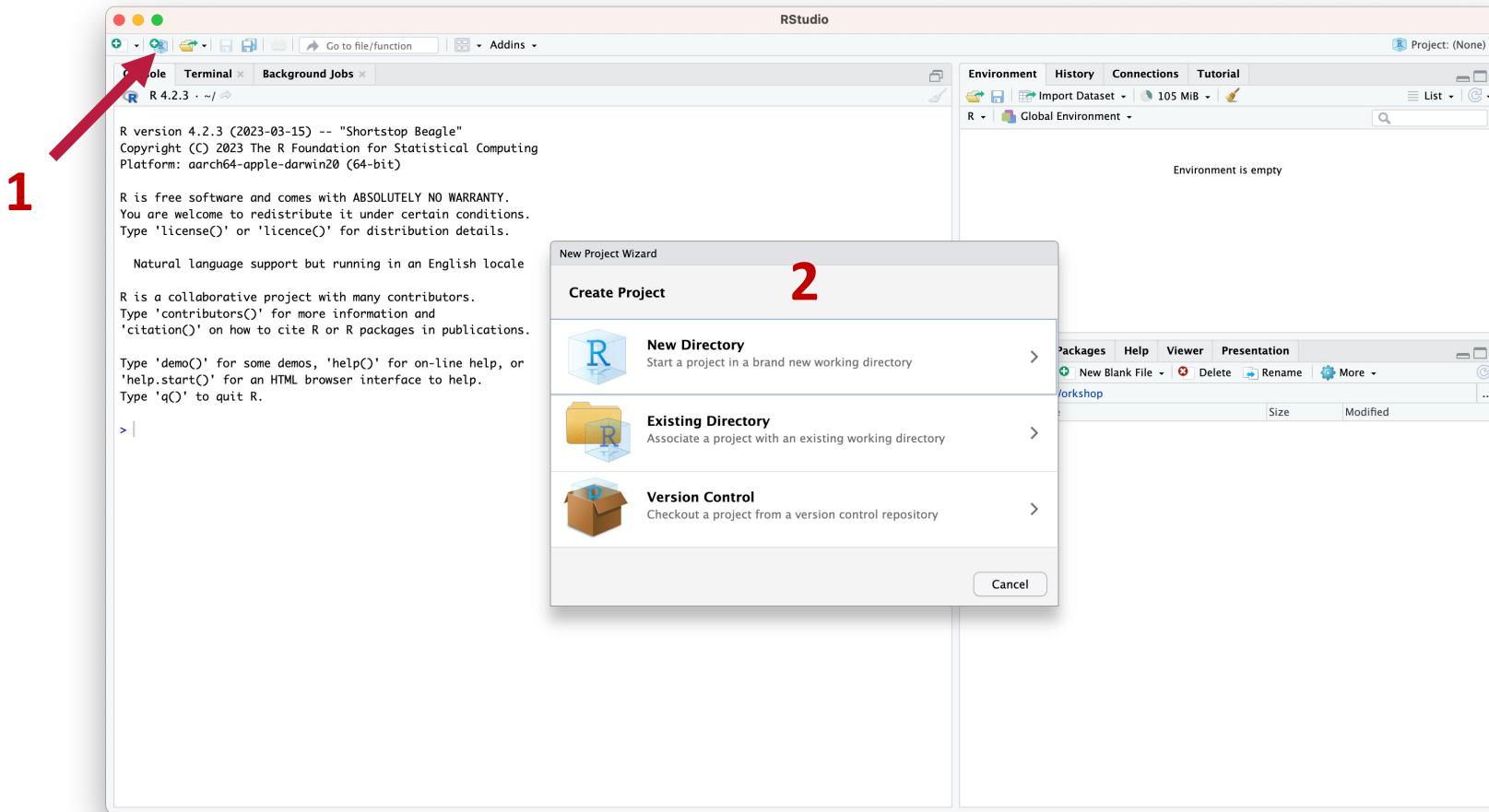
*Or just use `drop_na()` from the `tidyverse` package:

```
> head(airquality)
   Ozone Solar.R Wind Temp Month Day
1    41     190  7.4   67     5   1
2    36     118  8.0   72     5   2
3    12     149 12.6   74     5   3
4    18     313 11.5   62     5   4
5    NA      NA 14.3   56     5   5
6    28      NA 14.9   66     5   6
> tidyverse::drop_na(head(airquality))
   Ozone Solar.R Wind Temp Month Day
1    41     190  7.4   67     5   1
2    36     118  8.0   72     5   2
3    12     149 12.6   74     5   3
4    18     313 11.5   62     5   4
```

Project Time!

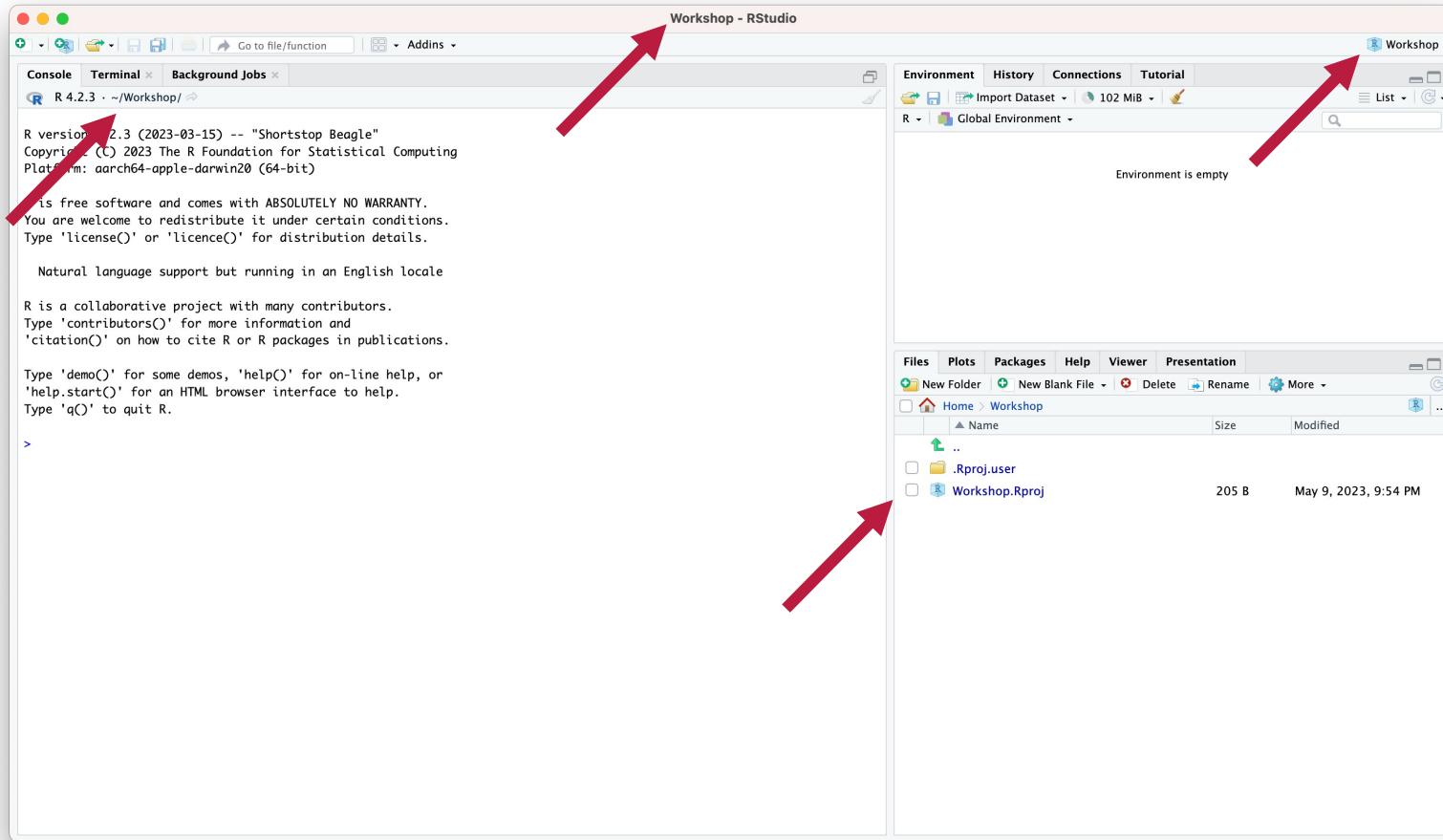
RStudio Projects

RStudio projects make it straightforward to divide your work into multiple contexts, each with their own working directory, workspace, history, and source documents.



RStudio Projects

RStudio projects make it straightforward to divide your work into multiple contexts, each with their own working directory, workspace, history, and source documents.



RStudio Projects

When a project is opened within RStudio the following actions are taken:

- A new R session (process) is started
- The `.Rprofile` file in the project's main directory (if any) is sourced by R
- The `.RData` file in the project's main directory is loaded (if project options indicate that it should be loaded).
- The `.Rhistory` file in the project's main directory is loaded into the RStudio History pane (and used for Console Up/Down arrow command history).
- The current working directory is set to the project directory.
- Previously edited source documents are restored into editor tabs
- Other RStudio settings (e.g., active tabs, splitter positions, etc.) are restored to where they were the last time the project was closed.

Reproducible Environments with the `renv` Package

Isolated

- Installing a new or updated package for one project won't break your other projects, and vice versa. That's because `renv` gives each project its own private package library.

Portable

- Easily transport your projects from one computer to another, even across different platforms. `renv` makes it easy to install the packages your project depends on.

Reproducible

- `renv` records the exact R and package versions you depend on, and ensures those exact versions are the ones that get installed wherever you go.

<https://docs.posit.co/ide/user/ide/guide/environments/r/renv.html>

renv General Workflow

1. Call `renv::init()` to initialize a new project-local environment with a private R library,
2. Work in the project as normal, installing and removing new R packages as they are needed in the project,
3. Call `renv::snapshot()` to save the state of the project library to the lockfile (called `renv.lock`),
4. Continue working on your project, installing and updating R packages as needed.
5. Call `renv::snapshot()` again to save the state of your project library if your attempts to update R packages were successful
6. Call `renv::restore()` to revert to the previous state as encoded in the lockfile if your attempts to update packages introduced some new problems.

<https://docs.posit.co/ide/user/ide/guide/environments/r/renv.html>

renv Cache Location

Platform	Location
Linux	<code>~/.cache/R/renv</code>
macOS	<code>~/Library/Caches/org.R-project.R/R/renv</code>
Windows	<code>%LOCALAPPDATA%/R/cache/R/renv</code>

Data Transformation with the Tidyverse

Toy Data

Let's look at our toy data environmental_data.csv

```
# A tibble: 124,744 × 6
# Groups:   date, time, barn [62,372]
  date      time    barn location     rh   temp
  <date>    <time> <chr>   <chr>   <dbl> <dbl>
1 2022-01-01 11'58" lactating inside     74  7.74
2 2022-01-01 11'58" lactating outside    87  2.5 
3 2022-01-01 11'58" sp_needs  inside     78 10.1 
4 2022-01-01 11'58" sp_needs  outside    87  2.5 
5 2022-01-01 26'58" lactating inside     74  8.31 
6 2022-01-01 26'58" lactating outside    87  2.5 
7 2022-01-01 26'58" sp_needs  inside     77  9.95 
8 2022-01-01 26'58" sp_needs  outside    87  2.5 
9 2022-01-01 41'58" lactating inside     74  8.89 
10 2022-01-01 41'58" lactating outside   87  2.5 
# ... with 124,734 more rows
# i Use `print(n = ...)` to see more rows
```

Ultimate Question

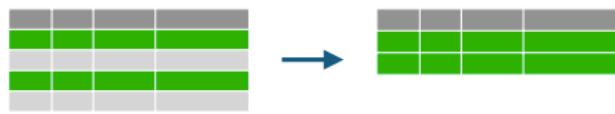
Are the temperatures inside the barns milder than outside?

i.e., warmer in the winter and colder in the summer

Isolating Data



Extract variables with **`select()`**



Extract observations with **`filter()`**



Arrange/Sort observations with **`arrange()`**

filter()

Extract rows that meet logical criteria

`filter(.data, ...)`

`.data` dataframe to transform

`...` one or more logical tests (filter returns each row for which the test is TRUE)

filter()

Extract rows that meet logical criteria

```
filter(env_data, barn == "lactating")
```

```
# A tibble: 124,744 × 6
# Groups: date, time, barn [62,372]
  date      time    barn location     rh   temp
  <date>    <time> <chr>  <chr>   <dbl> <dbl>
1 2022-01-01 11'58" lactating inside    74  7.74
2 2022-01-01 11'58" lactating outside   87  2.5
3 2022-01-01 11'58" sp_needs inside    78 10.1
4 2022-01-01 11'58" sp_needs outside   87  2.5
5 2022-01-01 26'58" lactating inside    74  8.31
6 2022-01-01 26'58" lactating outside   87  2.5
7 2022-01-01 26'58" sp_needs inside    77  9.95
8 2022-01-01 26'58" sp_needs outside   87  2.5
9 2022-01-01 41'58" lactating inside    74  8.89
10 2022-01-01 41'58" lactating outside  87  2.5
```



```
# A tibble: 62,372 × 6
  date      time    barn location     rh   temp
  <date>    <time> <chr>  <chr>   <dbl> <dbl>
1 2022-01-01 00:11:58 lactating inside    74  7.74
2 2022-01-01 00:11:58 lactating outside   87  2.5
3 2022-01-01 00:26:58 lactating inside    74  8.31
4 2022-01-01 00:26:58 lactating outside   87  2.5
5 2022-01-01 00:41:58 lactating inside    74  8.89
6 2022-01-01 00:41:58 lactating outside   87  2.5
7 2022-01-01 00:56:58 lactating inside    74  8.28
8 2022-01-01 00:56:58 lactating outside   87  2.6
9 2022-01-01 01:11:58 lactating inside    75  8.6
10 2022-01-01 01:11:58 lactating outside  87  2.8
```

filter()

Extract rows that meet logical criteria

```
filter(env_data, barn == "lactating")
```

== tests if equal
(returns TRUE or FALSE)

= used to set things
(returns nothing)

```
# A tibble: 124,744 × 6
# Groups: date, time, barn [62,372]
  date      time    barn location     rh   temp
  <date>    <time> <chr>  <chr>    <dbl> <dbl>
1 2022-01-01 11'58" lactating inside    74  7.74
2 2022-01-01 11'58" lactating outside   87  2.5
3 2022-01-01 11'58" sp_needs inside    78 10.1
4 2022-01-01 11'58" sp_needs outside   87  2.5
5 2022-01-01 26'58" lactating inside    74  8.31
6 2022-01-01 26'58" lactating outside   87  2.5
7 2022-01-01 26'58" sp_needs inside    77  9.95
8 2022-01-01 26'58" sp_needs outside   87  2.5
9 2022-01-01 41'58" lactating inside    74  8.89
10 2022-01-01 41'58" lactating outside  87  2.5
```



```
# A tibble: 62,372 × 6
  date      time    barn location     rh   temp
  <date>    <time> <chr>  <chr>    <dbl> <dbl>
1 2022-01-01 00:11:58 lactating inside    74  7.74
2 2022-01-01 00:11:58 lactating outside   87  2.5
3 2022-01-01 00:26:58 lactating inside    74  8.31
4 2022-01-01 00:26:58 lactating outside   87  2.5
5 2022-01-01 00:41:58 lactating inside    74  8.89
6 2022-01-01 00:41:58 lactating outside   87  2.5
7 2022-01-01 00:56:58 lactating inside    74  8.28
8 2022-01-01 00:56:58 lactating outside   87  2.6
9 2022-01-01 01:11:58 lactating inside    75  8.6
10 2022-01-01 01:11:58 lactating outside  87  2.8
```

Logical Tests

$x < y$	Less than
$x > y$	Greater than
$x == y$	Equals to
$x <= y$	Less than or equal to
$x >= y$	Greater than or equal to
$x != y$	Not equal to

$x \%in\% y$	x is at least one of y
<code>is.na(x)</code>	Is NA
<code>!is.na(x)</code>	Is not NA
$a \& b$	and
$a b$	or

Your Turn!

Use the logical operators to manipulate our env_data to show observations that are:

- From inside the barns
- Above 30°C
- Between 0°C and 10°C inside the Special Needs barn
- From June 3, 2022, or June 4, 2022

Your Turn!

Use the logical operators to manipulate our env_data to show observations that are:

```
### From inside the barn ----  
filter(env_data, location == "inside")  
  
### Above 30C ----  
filter(env_data, temp > 30)  
  
### Between 0 and 10C inside the Special Needs barn ----  
filter(env_data, temp >= 0, temp <= 10, location == "inside",  
      barn == "sp_needs")  
  
### From June 3, 2022, or June 4, 2022 ----  
filter(env_data, date == "2022-06-04" | date == "2022-06-05")
```

select()

Extract columns by name

```
select(.data, ...)
```

.data dataframe to transform

... name(s) of columns to extract, or a select helper function

select()

Only show records for temperature

```
select(env_data, date, time, barn, location, temp)
```

```
# A tibble: 124,744 × 5
  date      time    barn   location  temp
  <date>    <time>  <chr>   <chr>     <dbl>
1 2022-01-01 11'58" lactating inside    7.74
2 2022-01-01 11'58" lactating outside   2.5
3 2022-01-01 11'58" sp_needs  inside   10.1
4 2022-01-01 11'58" sp_needs  outside   2.5
5 2022-01-01 26'58" lactating inside    8.31
6 2022-01-01 26'58" lactating outside   2.5
7 2022-01-01 26'58" sp_needs  inside   9.95
8 2022-01-01 26'58" sp_needs  outside   2.5
9 2022-01-01 41'58" lactating inside    8.89
10 2022-01-01 41'58" lactating outside   2.5
```

select() helpers

- Select a range of columns (:)
 - `select(env_data, date:location)`
- Select every column but (-)
 - `select(env_data, -rh)`
- Select columns that start with something (`starts_with()`)
 - `select(env_data, starts_with("t"))`
- Select columns that end with something (`ends_with()`)
 - `select(env_data, ends_with("e"))`

select() helpers

And a few more!

Data transformation with dplyr :: CHEAT SHEET

dplyr functions work with pipes and expect **tidy data**. In tidy data:
Each variable is in its own column & Each observation, or case, is in its own row x %>% f(y) becomes f(x, y)

Summarise Cases

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).

summary function

summarise(data, ...)
Compute table of summaries.
summarise(mtcars, avg = mean(mpg))

count(data, ..., wt = NULL, sort = FALSE, name = NULL) Count the number of rows in each group defined by the variables in ... Also **tally**.

Group Cases

Use **group_by**(data, ..., add = FALSE, drop = TRUE) to create a "grouped" copy of a table grouped by columns. In dplyr functions will manipulate each "group" separately and combine the results.

mtcars %>%
group_by(cyl) %>%
summarise(avg = mean(mpg))

Use **rowwise**(data, ...) to group data into individual rows. dplyr functions will compute results for each row. Also apply functions to list-columns. See tidy cheat sheet for list-column workflow.

starwars %>%
rowwise() %>%
mutate(lm_count = length(films))

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



Manipulate Cases

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

filter(data, ..., preserve = FALSE) Extract rows that meet logical criteria.
filter(mtcars, mpg > 28)
distinct(data, ..., keep_all = FALSE) Remove rows with duplicate values.
distinct(mtcars, gear)
slice(data, ..., preserve = FALSE) Select rows by position.
slice(mtcars, 10:15)
slice_sample(data, ..., n, prop, weight_by = NULL, replace = FALSE) Randomly select rows. Use n to select a number of rows and prop to select a fraction of rows.
slice_sample(mtcars, n = 5, replace = TRUE)
slice_min(data, order_by, ..., n, prop, with_ties = TRUE) and slice_max() Select rows with the lowest and highest values.
slice_min(mtcars, mpg, prop = 0.25)
slice_head(data, ..., n, prop) and slice_tail() Select the first or last rows.
slice_head(mtcars, n = 5)

Logical and boolean operators to use with filter()
== < <= > >= is.na() %in% | xor()
!= > >= !is.na() ! &
See ?base::Logic and ?Comparison for help.

ARRANGE CASES
arrange(data, ..., by = group = FALSE) Order rows by values of a column or columns (low to high), use with **desc**() to order from high to low.
arrange(mtcars, mpg)
arrange(mtcars, desc(mpg))

ADD CASES
add_row(data, ..., before = NULL, .after = NULL) Add one or more rows to a table.
add_row(cars, speed = 1, dist = 1)

Manipulate Variables

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

pull(data, var = -1, name = NULL, ...) Extract column values as a vector, by name or index.
pull(mtcars, wt)
select(data, ...) Extract columns as a table.
select(mtcars, mpg, wt)
relocate(data, ..., before = NULL, .after = NULL) Move columns to new position.
relocate(mtcars, mpg, cyl, .after = last_col())

Use these helpers with **select()** and **across()**
e.g. select(mtcars, mpg:cyl)
contains(match) num_range(prefix, range) : e.g. mpg:cyl
ends_with(match) all_of(x)/any_of(x, ..., vars) - e.g. -gear
starts_with(match) matches(match)

MANIPULATE MULTIPLE VARIABLES AT ONCE
across(cols, funs, ..., .names = NULL) Summarise or mutate multiple columns in the same way.
summarise(mtcars, across(everything(), mean))
c_across(cols) Compute across columns in row-wise data.
transmute(rowwise(UKgas), total = sum(c_across(1:2)))

MAKE NEW VARIABLES
Apply **vectorized functions** to columns. Vectorized functions take vectors as input and return vectors of the same length as output (see back).
vectorized function

arrange(data, ..., by = group = FALSE) Order rows by values of a column or columns (low to high), use with **desc**() to order from high to low.
arrange(mtcars, mpg)
arrange(mtcars, desc(mpg))

mutate(data, ..., keep = "all", .before = NULL, .after = NULL) Compute new column(s). Also **add**, **create**, **add_count**, and **add_tally**().
mutate(mtcars, gpm = 1 / mpg)

transmute(data, ...) Compute new column(s), drop others.
transmute(mtcars, gpm = 1 / mpg)

rename(data, ...) Rename columns. Use **rename** with a function.
rename(cars, distance = dist)

EXTRACT VARIABLES

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

pull(.data, var = -1, name = NULL, ...) Extract column values as a vector, by name or index.
pull(mtcars, wt)

select(.data, ...) Extract columns as a table.
select(mtcars, mpg, wt)

relocate(.data, ..., .before = NULL, .after = NULL)
Move columns to new position.
relocate(mtcars, mpg, cyl, .after = last_col())

Use these helpers with **select() and **across()****
e.g. select(mtcars, mpg:cyl)

contains(match) **num_range(prefix, range)** : e.g. mpg:cyl
ends_with(match) **all_of(x)/any_of(x, ..., vars)** - e.g. -gear
starts_with(match) **matches(match)**

Your Turn!

Use filter() and select() to show only relative humidity from inside the Special Needs barn

*hint: feel free to create an intermediate object

```
### Relative humidity from inside the Special Needs barn ----  
sp_needs <- filter(env_data, barn == "sp_needs", location == "inside")  
sp_needs <- select(sp_needs, date, time, rh)  
sp_needs  
  
# A tibble: 31,186 × 3  
  date        time       rh  
  <date>     <time>     <dbl>  
1 2022-01-01 00:11:58    78  
2 2022-01-01 00:26:58    77  
3 2022-01-01 00:41:58    77  
4 2022-01-01 00:56:58    78  
5 2022-01-01 01:11:58    78
```

arrange()

Order rows from smallest to largest values

`arrange(.data, ...)`

`.data` dataframe to transform

`...` One or more columns to order by (additional columns will be used as tie breakers)

arrange()

1. Order by **temp**

```
arrange(env_data, temp)
```

```
# A tibble: 124,744 × 6
  date      time    barn   location     rh   temp
  <date>    <time>  <chr>  <chr>     <dbl> <dbl>
1 2022-01-21 08:11:58 sp_needs outside  77 -24
2 2022-01-21 08:11:58 lactating outside 77 -23.9
3 2022-01-29 05:41:58 sp_needs outside  71 -23.9
4 2022-01-29 05:41:58 lactating outside 71 -23.8
5 2022-01-29 07:26:58 sp_needs outside  72 -23.8
6 2022-01-29 07:56:58 sp_needs outside  73 -23.8
7 2022-01-29 07:56:58 lactating outside 73 -23.7
8 2022-01-21 05:41:58 lactating outside 76 -23.6
9 2022-01-21 05:41:58 sp_needs outside 77 -23.6
10 2022-01-29 07:26:58 lactating outside 72 -23.6
```

2. Order by **temp** descending

```
arrange(env_data, desc(temp))
```

```
# A tibble: 124,744 × 6
  date      time    barn   location     rh   temp
  <date>    <time>  <chr>  <chr>     <dbl> <dbl>
1 2022-09-08 15:11:58 sp_needs outside  18  39.2
2 2022-06-24 15:41:58 lactating outside 25  39.1
3 2022-06-24 15:41:58 sp_needs outside 25  39.1
4 2022-08-15 14:41:58 lactating outside 18  39.1
5 2022-08-15 14:41:58 sp_needs outside 18  39.1
6 2022-09-08 15:11:58 lactating outside 18  39.1
7 2022-06-24 15:56:58 sp_needs outside 24  39
8 2022-06-24 15:56:58 lactating outside 24  38.9
9 2022-09-08 15:41:58 sp_needs outside 18  38.9
10 2022-09-08 15:56:58 lactating outside 18  38.9
```

Your Turn!

Order by **temp** and use **rh** as tie breaker. What was the lowest temperature?

Your Turn!

Order by **temp** and use **rh** as tie breaker. What was the lowest temperature?

```
### Order by temp and use rh as tie breaker ----  
arrange(env_data, temp, rh)
```

```
# A tibble: 124,744 × 6  
  date      time    barn location   rh   temp  
  <date>    <time>  <chr>  <chr>  <dbl> <dbl>  
1 2022-01-21 08:11:58 sp_needs outside  77 -24  
2 2022-01-21 08:11:58 lactating outside 77 -23.9  
3 2022-01-29 05:41:58 sp_needs outside  71 -23.9  
4 2022-01-29 05:41:58 lactating outside 71 -23.8  
5 2022-01-29 07:26:58 sp_needs outside  72 -23.8  
6 2022-01-29 07:56:58 sp_needs outside  73 -23.8  
7 2022-01-29 07:56:58 lactating outside 73 -23.7  
8 2022-01-21 05:41:58 lactating outside 76 -23.6  
9 2022-01-21 05:41:58 sp_needs outside 77 -23.6  
10 2022-01-29 07:26:58 lactating outside 72 -23.6
```

```
# A tibble: 124,744 × 6  
  date      time    barn location   rh   temp  
  <date>    <time>  <chr>  <chr>  <dbl> <dbl>  
1 2022-01-21 08:11:58 sp_needs outside  77 -24  
2 2022-01-29 05:41:58 sp_needs outside 71 -23.9  
3 2022-01-21 08:11:58 lactating outside 77 -23.9  
4 2022-01-29 05:41:58 lactating outside 71 -23.8  
5 2022-01-29 07:26:58 sp_needs outside 72 -23.8  
6 2022-01-29 07:56:58 sp_needs outside 73 -23.8  
7 2022-01-29 07:56:58 lactating outside 73 -23.7  
8 2022-01-29 07:26:58 lactating outside 72 -23.6  
9 2022-01-29 07:41:58 sp_needs outside 73 -23.6  
10 2022-01-21 05:41:58 lactating outside 76 -23.6
```

The Pipe Operator %>%

Use filter(), select() and arrange() to show only relative humidity from inside the Special Needs barn, ordered ascending by relative humidity

```
sp_needs <- filter(env_data, barn == "sp_needs", location == "inside")
sp_needs <- select(sp_needs , date, time, rh)
sp_needs <- arrange(sp_needs, rh)
sp_needs
```

```
# A tibble: 31,186 × 3
  date      time      rh
  <date>    <time>    <dbl>
1 2022-05-12 13:11:58 29
2 2022-05-12 15:26:58 29
3 2022-05-12 16:11:58 29
4 2022-05-12 12:56:58 30
5 2022-05-12 13:26:58 30
6 2022-05-12 13:41:58 30
7 2022-05-12 13:56:58 30
8 2022-05-12 14:11:58 30
9 2022-05-12 14:26:58 30
10 2022-05-12 15:56:58 30
```

The Pipe Operator %>%

Use filter(), select() and arrange() to show only relative humidity from inside the Special Needs barn, ordered ascending by relative humidity

```
sp_needs <- filter(env_data, barn == "sp_needs", location == "inside")  
sp_needs <- select(sp_needs , date, time, rh)  
sp_needs <- arrange(sp_needs, rh)  
sp_needs
```

	date	time	rh
	<date>	<time>	<dbl>
1	2022-05-12	13:11:58	29
2	2022-05-12	15:26:58	29
3	2022-05-12	16:11:58	29
4	2022-05-12	12:56:58	30
5	2022-05-12	13:26:58	30
6	2022-05-12	13:41:58	30
7	2022-05-12	13:56:58	30
8	2022-05-12	14:11:58	30
9	2022-05-12	14:26:58	30
10	2022-05-12	15:56:58	30

The Pipe Operator %>%

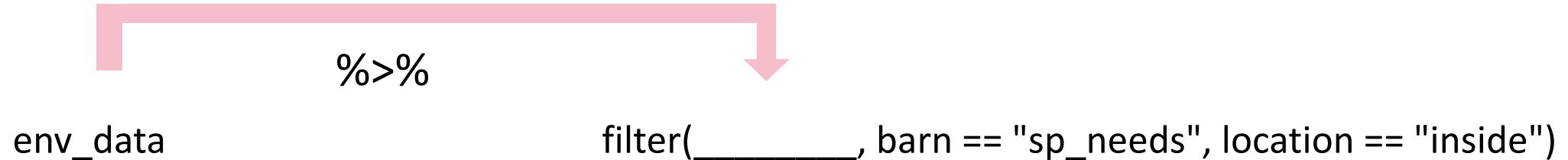
Use filter(), select() and arrange() to show only relative humidity from inside the Special Needs barn, ordered ascending by relative humidity

```
sp_needs <- filter(env_data, barn == "sp_needs", location == "inside")
sp_needs <- select(sp_needs , date, time, rh)
sp_needs <- arrange(sp_needs, rh)
sp_needs
```

```
sp_needs <- arrange(select(filter(env_data, barn == "sp_needs", location ==
"inside"), date, time, rh), rh)
```

The Pipe Operator %>%

The pipe operator %>% passes result on the left into first argument of function on the right.



The Pipe Operator %>%

The pipe operator %>% passes result on the left into first argument of function on the right.



```
env_data %>% filter(_____, barn == "sp_needs", location == "inside")
```

```
filter(env_data, barn == "sp_needs", location == "inside")
```

equals to

```
env_data %>% filter(barn == "sp_needs", location == "inside")
```

Your Turn!

Use filter(), select(), arrange() and the pipe operator %>% to show only relative humidity from inside the Special Needs barn, ordered ascending by relative humidity

```
env_data %>%
  filter(barn == "sp_needs", location == "inside") %>%
  select(date, time, rh) %>%
  arrange(rh)
# A tibble: 31,186 × 3
  date       time     rh
  <date>     <time>   <dbl>
1 2022-05-12 13:11:58 29
2 2022-05-12 15:26:58 29
3 2022-05-12 16:11:58 29
4 2022-05-12 12:56:58 30
5 2022-05-12 13:26:58 30
6 2022-05-12 13:41:58 30
7 2022-05-12 13:56:58 30
8 2022-05-12 14:11:58 30
9 2022-05-12 14:26:58 30
10 2022-05-12 15:56:58 30
```

Ultimate Question

Are the temperatures inside the barns milder than outside?

i.e., warmer in the winter and colder in the summer

Ultimate Question

Are the temperatures inside the barns milder than outside?
i.e., warmer in the winter and colder in the summer

What do we need to know?

- Average temperatures during winter and summer months for each barn, inside and outside

barn	season	location	avg_temp
sp_needs	winter	inside	
sp_needs	winter	outside	
sp_needs	summer	inside	
sp_needs	summer	outside	
lactating	winter	inside	
lactating	winter	outside	
lactating	summer	inside	
lactating	summer	outside	

Ultimate Question

Are the temperatures inside the barns milder than outside?
i.e., warmer in the winter and colder in the summer

What do we need to know?

- Average temperatures during winter and summer months for each barn, inside and outside

barn	season	location	avg_temp
sp_needs	winter	inside	
sp_needs	winter	outside	
sp_needs	summer	inside	
sp_needs	summer	outside	
lactating	winter	inside	
lactating	winter	outside	
lactating	summer	inside	
lactating	summer	outside	

Deriving Information

We can make table of summaries with **summarize()/summarise()**

We can make new variables with **mutate()**

summarize()

Transforms a vector of data into one value

```
summarize(.data, new_column = function(vector))
```

.data	dataframe to transform
new_column	New column created by function()
function(...)	Function used to transform a vector
vector	Vector to be transformed, it's a column from .data

summarize()

Use summarize() to create a summary table with average and max temperatures

```
env_data %>%  
  summarize(avg_temp = mean(temp),  
           max_temp = max(temp))
```

A tibble: 1 × 2
 avg_temp max_temp
 <dbl> <dbl>
1 NA NA

```
env_data %>% filter(is.na(temp))
```

A tibble: 8 × 6
 date time barn location rh temp
 <date> <time> <chr> <chr> <dbl> <dbl>
1 2022-01-31 10:41:58 lactating inside NA NA
2 2022-01-31 10:41:58 lactating outside NA NA
3 2022-01-31 10:41:58 sp_needs inside NA NA
4 2022-01-31 10:41:58 sp_needs outside NA NA
5 2022-01-31 10:56:58 lactating inside NA NA
6 2022-01-31 10:56:58 lactating outside NA NA
7 2022-01-31 10:56:58 sp_needs inside NA NA
8 2022-01-31 10:56:58 sp_needs outside NA NA

summarize()

Use summarize() to create a summary table with average and max temperatures

```
env_data %>%  
  summarize(avg_temp = mean(temp, na.rm = TRUE),  
            max_temp = max(temp, na.rm = TRUE))
```

OR

```
env_data %>%  
  na.omit() %>%  
  summarize(avg_temp = mean(temp),  
            max_temp = max(temp))
```

```
# A tibble: 1 × 2  
  avg_temp max_temp  
      <dbl>     <dbl>  
1       .      12.2     39.2
```

Your Turn!

Using summarize() and filter(), get the min, mean, and max temperatures and relative humidity inside the Special Needs barn.

* Don't forget to account for NAs

```
env_data %>%
  na.omit() %>%
  filter(barn == "sp_needs",
         location == "inside") %>%
  summarize(min_temp = min(temp),
            avg_temp = mean(temp),
            max_temp = max(temp),
            min_rh = min(rh),
            avg_rh = mean(rh),
            max_rh = max(rh))
```

A tibble: 1 × 6

	min_temp	avg_temp	max_temp	min_rh	avg_rh	max_rh
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1	-1.95	14.4	31.4	29	73.3	91

Ultimate Question

Are the temperatures inside the barns milder than outside?
i.e., warmer in the winter and colder in the summer

What do we need to know?

- Average temperatures during winter and summer months for each barn, inside and outside

barn	season	location	avg_temp
sp_needs	winter	inside	
sp_needs	winter	outside	
sp_needs	summer	inside	
sp_needs	summer	outside	
lactating	winter	inside	
lactating	winter	outside	
lactating	summer	inside	
lactating	summer	outside	

group_by()

Group observations by common values

```
group_by(.data, ...)
```

.data dataframe to transform

... one or more column names to group

group_by()

Get the min, mean, and max temperatures per barn and location

```
env_data %>%  
  na.omit() %>%  
  summarize(min_temp = min(temp),  
            avg_temp = mean(temp),  
            max_temp = max(temp))  
  
# A tibble: 1 × 3  
#>   min_temp avg_temp max_temp  
#>   <dbl>     <dbl>     <dbl>  
#> 1      -24     12.2     39.2
```

group_by()

Get the min, mean, and max temperatures per barn and location

```
env_data %>%  
  na.omit() %>%  
  group_by(barn, location) %>%  
  summarize(min_temp = min(temp),  
            avg_temp = mean(temp),  
            max_temp = max(temp))
```

# A tibble: 4 × 5					
# Groups: barn [2]					
barn	location	min_temp	avg_temp	max_temp	
<chr>	<chr>	<dbl>	<dbl>	<dbl>	
1	lactating inside	1.34	13.8	31.3	
2	lactating outside	-23.9	10.3	39.1	
3	sp_needs inside	-1.95	14.4	31.4	
4	sp_needs outside	-24	10.3	39.2	

group_by()

```
sampled_env_data %>% na.omit() %>%  
  summarize(min = min(temp), avg = mean(temp), max = max(temp))
```

barn	location	rh	temp
lactating	inside	79	20.4
lactating	inside	85	19.9
lactating	outside	83	-2.5
lactating	outside	47	23.4
sp_needs	inside	83	12.0
sp_needs	inside	78	9.2
sp_needs	outside	73	-14.5
sp_needs	outside	58	15.7



min	avg	max
-14.5	10.4	23.4

group_by()

group_by() + summarize()

barn	location	rh	temp
lactating	inside	79	20.4
lactating	inside	85	19.9
lactating	outside	83	-2.5
lactating	outside	47	23.4
sp_needs	inside	83	12.0
sp_needs	inside	78	9.2
sp_needs	outside	73	-14.5
sp_needs	outside	58	15.7



min	avg	max
19.9	20.1	20.4



-2.5	10.4	23.4
------	------	------



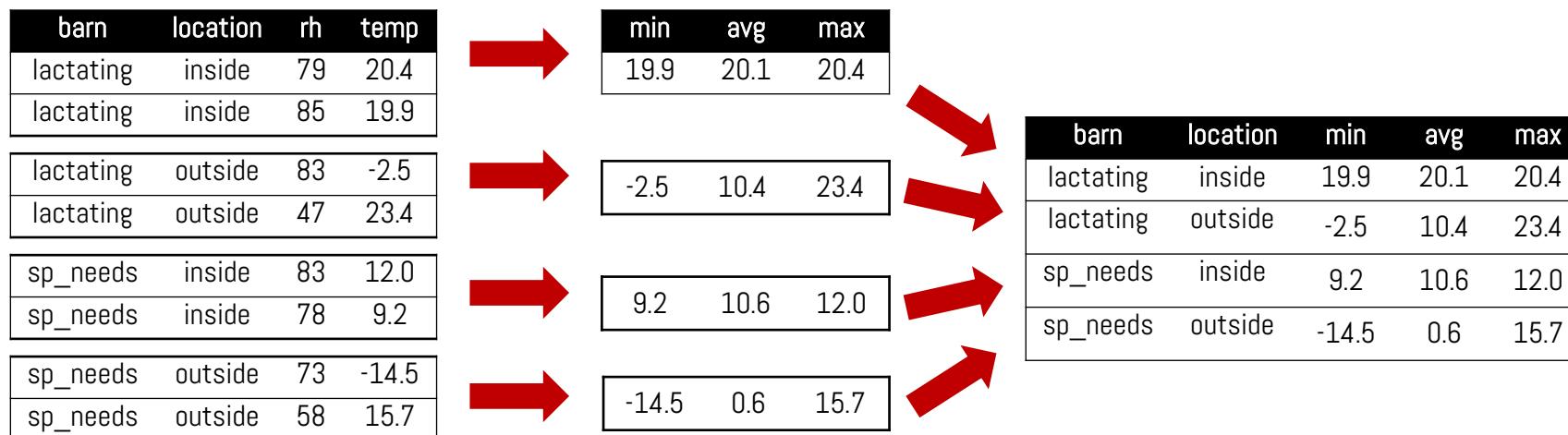
9.2	10.6	12.0
-----	------	------



-14.5	0.6	15.7
-------	-----	------

group_by()

```
sampled_env_data %>% group_by(barn, location) %>% na.omit() %>%  
  summarize(min = min(temp),  
            avg = mean(temp),  
            max = max(temp))
```



Your Turn!

Use group_by(), filter(), and summarize() to show the lowest and highest relative humidity and temperature of the inside of each barn

```
env_data %>%  
  filter(location == "inside") %>%  
  na.omit() %>%  
  group_by(barn, location) %>%  
  summarize(min_temp = min(temp),  
            max_temp = max(temp),  
            min_rh = min(rh),  
            max_rh = max(rh))
```

```
A tibble: 2 × 6  
Groups: barn [2]  
#> barn     location min_temp max_temp min_rh max_rh  
#> <chr>    <chr>      <dbl>     <dbl>    <dbl>   <dbl>  
#> lactating inside      1.34      31.3     24     90  
#> sp_needs  inside     -1.95      31.4     29     91
```

Ultimate Question

Are the temperatures inside the barns milder than outside?
i.e., warmer in the winter and colder in the summer

What do we need to know?

- Average temperatures during winter and summer months for each barn, inside and outside

barn	season	location	avg_temp
sp_needs	winter	inside	
sp_needs	winter	outside	
sp_needs	summer	inside	
sp_needs	summer	outside	
lactating	winter	inside	
lactating	winter	outside	
lactating	summer	inside	
lactating	summer	outside	

mutate()

Apply vectorized functions to columns to create new columns

```
mutate(.data, new_column = function(vector))
```

.data dataframe to transform

new_column New column created by function()

function(...) Function used to transform a vector

vector Vector to be transformed, can be a column from .data

mutate()

Create new columns

```
env_data %>%  
  mutate(year = lubridate::year(date),  
         month = lubridate::month(date),  
         day = lubridate::day(date),  
         barn = dplyr::if_else(barn == "sp_needs", "special_needs", barn))
```

```
# A tibble: 124,744 × 9  
  date      time    barn     location   rh  temp  year month   day  
  <date>    <time>  <chr>    <chr>    <dbl> <dbl> <dbl> <dbl> <int>  
1 2022-01-01 11'58" lactating  inside     74  7.74 2022     1     1  
2 2022-01-01 11'58" lactating  outside    87  2.5  2022     1     1  
3 2022-01-01 11'58" special_needs inside    78 10.1 2022     1     1  
4 2022-01-01 11'58" special_needs outside   87  2.5  2022     1     1  
5 2022-01-01 26'58" lactating  inside     74  8.31 2022     1     1
```

Your Turn!

Using the functions from today, create the following table to answer our ultimate question:

Are the temperatures inside the barns milder than outside?

barn	season	location	avg_temp
sp_needs	winter	inside	
sp_needs	winter	outside	
sp_needs	summer	inside	
sp_needs	summer	outside	
lactating	winter	inside	
lactating	winter	outside	
lactating	summer	inside	
lactating	summer	outside	

Winter: Between 2021-12-31 & 2022-03-20
Summer: Between 2022-06-21 & 2022-09-23
Spring/fall: Anything else

Your Turn!

Using the functions from today, create the following table to answer our ultimate question:

Are the temperatures inside the barns milder than outside?

```
env_data %>%
  na.omit() %>%
  mutate(season = if_else(date >= "2021-12-31" & date <= "2022-03-20",
                         true = "winter",
                         false = if_else(date >= "2022-06-21" &
                                         date <= "2022-09-23",
                                         true = "summer",
                                         false = "spring/fall")) %>%
  filter(season %in% c("summer", "winter")) %>%
  group_by(barn, season, location) %>%
  summarise(avg_temp = mean(temp)) %>%
  arrange(desc(barn), desc(season))
```

	barn	season	location	avg_temp
	<chr>	<chr>	<chr>	<dbl>
1	sp_needs	winter	inside	7.87
2	sp_needs	winter	outside	-5.90
3	sp_needs	summer	inside	20.5
4	sp_needs	summer	outside	21.9
5	lactating	winter	inside	7.62
6	lactating	winter	outside	-5.88
7	lactating	summer	inside	20.3
8	lactating	summer	outside	21.8

YES!!

THANK YOU!

Questions?

Feel free to reach me later at
alcantal@uoguelph.ca



**AGRI-FOOD DATA
CANADA**

AT THE UNIVERSITY *of* GUELPH