

Introduction to R for Biologists



Gita Yadav

Introductions

Tutors

Participants Intro Document

Links- Course HandBook

Login to the OTE?

This course

- Aimed for **beginners** who never used R before.
 - Basics of R: Not Everything!
 - Tricks
 - How to get started
 - How you can learn more and improve yourself
 - Where to Find Help
- After this course you should be able to
 - Manipulate data and visualise it
 - Further your knowledge of R independently.
- As always! You need to practice frequently to get better in R

Time Table

Day 1	Day 2
Getting started	
Introduction to R	Data manipulation and visualisation with tidyverse
Starting with data	

RECOMMENDATION:

Use TWO screens to participate in this course!

One to Watch us

One for DIY!

SETUP INSTRUCTIONS

All the software required to run this course is installed on your OTE, so you do need to do anything.

In case you are not using the OTE, follow the instructions in COURSE HANDBOOK!

- R and RStudio are separate downloads and installations.
- R is the underlying statistical computing environment, but using R alone is no fun.
- RStudio is a graphical integrated development environment (IDE) that makes using R much easier and more interactive.
- You need to install R before you install RStudio.
- After installing both programs, you will need to install the **tidyverse** package from within RStudio.

SETUP INSTRUCTIONS

Dear Participant,

To take part in the exercise sessions during this R course you'll need to log into our training environment using the remote desktop interface provided. Click on the links below and enter the credentials given here to start. We'll go into further details during the course.

Please remember that we advise using **google chrome** and in particular avoiding the Safari browser for these exercises.

username='USERID' password= PASSWORD

[Remote Desktop Link](#)



CLICK THIS LINK!



UNIVERSITY OF
CAMBRIDGE

SETUP INSTRUCTIONS

Sign in to RStudio

Username:

Password:

Stay signed in when browser closes

You will automatically be signed out after 60 minutes of inactivity.

Sign In

SETUP INSTRUCTIONS

The screenshot shows the RStudio interface. The top bar includes tabs for Roman Empire, Discover, Google, PUBMED, NIPGR, Mail, Banking, RNA-Seq, Programs, RECIPES, SAMMY, Databases, Lectures, and Other Bookmarks. The main window has tabs for Console, Terminal, and Background Jobs. The Console tab displays the R startup message:

```
R version 4.1.2 (2021-11-01) -- "Bird Hippie"
Copyright (C) 2021 The R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

Session restored from your saved work on 2022-Aug-24 19:29:38 UTC (12 hours ago)
> |
```

The Environment tab shows the Global Environment is empty. The Files tab displays the following directory structure:

Name	Size	Modified
.r	0 B	Aug 23, 2022, 9:53 AM
.Rhistory		
Desktop		
Documents		
Downloads		
Music		
Pictures		
Public		
R		
snap		
Templates		
thinclient_drives		
Videos		





R is a statistical programming language.

It is very popular in the Data Science field, including Bioinformatics.

The term “R” is used to refer to both the programming language and the software that interprets the scripts written using



UNIVERSITY OF
CAMBRIDGE

Course Materials and Links

COURSE DATA

> Ecology. 2016 Apr;97(4):1082. doi: 10.1890/15-2115.1.

Long-term monitoring and experimental manipulation of a Chihuahuan desert ecosystem near Portal, Arizona (1977–2013)

S K Morgan Ernest ^{1 2}, Glenda M Yenni ^{1 2}, Ginger Allington ³, Erica M Christensen ^{1 2},
Keith Geluso ⁴, Jacob R Goheen ⁵, Michelle R Schutzenhofer ⁶, Sarah R Supp ⁷,
Katherine M Thibault ⁸, James H Brown ⁹, Thomas J Valone ¹⁰

Affiliations + expand

PMID: 28792597 DOI: 10.1890/15-2115.1

[Free article](#)

- Data is from the paper

[S. K. Morgan Ernest, Thomas J. Valone, and James H. Brown. 2009. Long-term monitoring and experimental manipulation of a Chihuahuan Desert ecosystem near Portal, Arizona, USA. Ecology 90:1708.](#)

- A simplified version of this data will be used in this course:

<https://doi.org/10.6084/m9.figshare.1314459>

We will download this directly from R during the lessons when we need them.

This course has been created by adapting and adding new content to the [Data Carpentry lesson in Ecology: 10.5281/zenodo.3264888](#)



UNIVERSITY OF
CAMBRIDGE



Why learn R?

- Run an R script with **one click/command.**



VS



Rscript <name of file>

[This Photo](#) by Unknown Author is licensed under [CC BY-NC](#)



UNIVERSITY OF
CAMBRIDGE



Why learn R?

- R promotes **Reproducibility**

Reproducibility is when someone else (including your future self) can obtain the same results from the same dataset when using the same analysis.

- Automating your analysis.
- Generate reports.



Why learn R?

- R is **interdisciplinary & Extensible!**
 - Over 10,000 packages that cover different fields
 - Bioconductor repository for bioinformatics packages (GENOMICS)
 - Image Analysis
 - GIS
 - Time Series
 - Population genetics
 - Machine Learning
 - and a lot more!



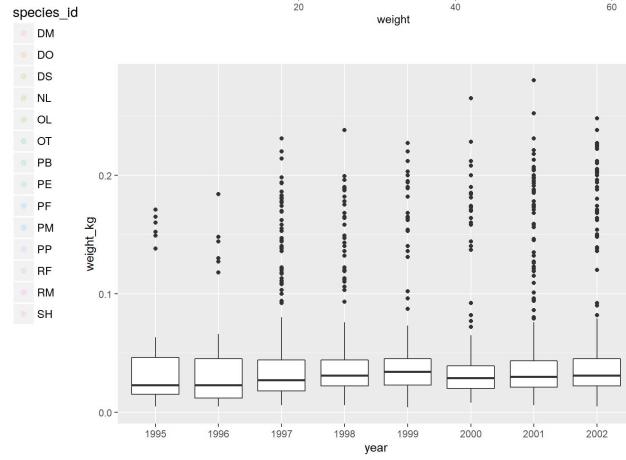
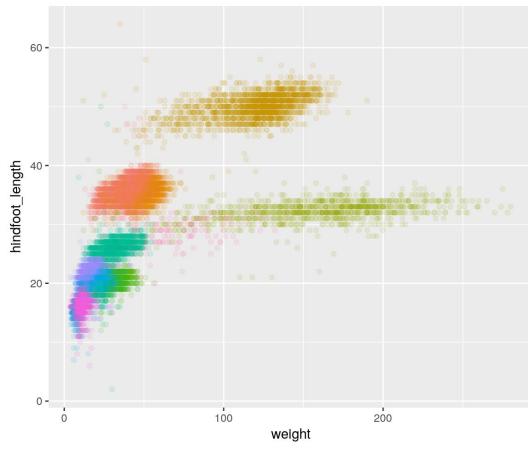
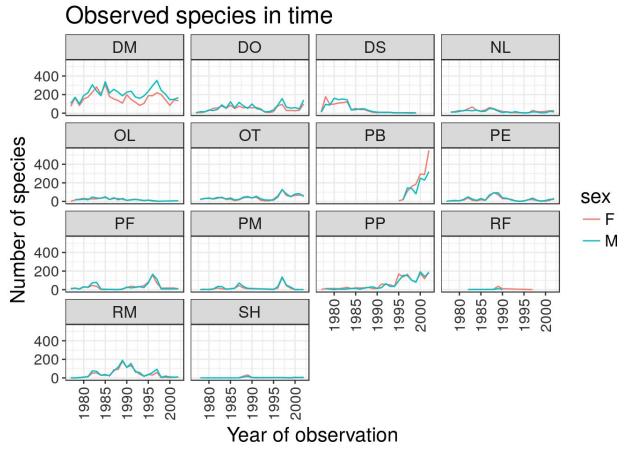
Why learn R?

- R works on data of different sizes (100s to Millions of Lines!)
- The skills you learn with R scale easily with the size of your dataset.
- Designed for data analysis: Special data structures and data types that make handling of missing data and statistical factors convenient.
- Seamless connection to spreadsheets, databases, and many other data formats, on your **computer** or on the **web**.



Why learn R?

- Best Feature: High-quality graphics!

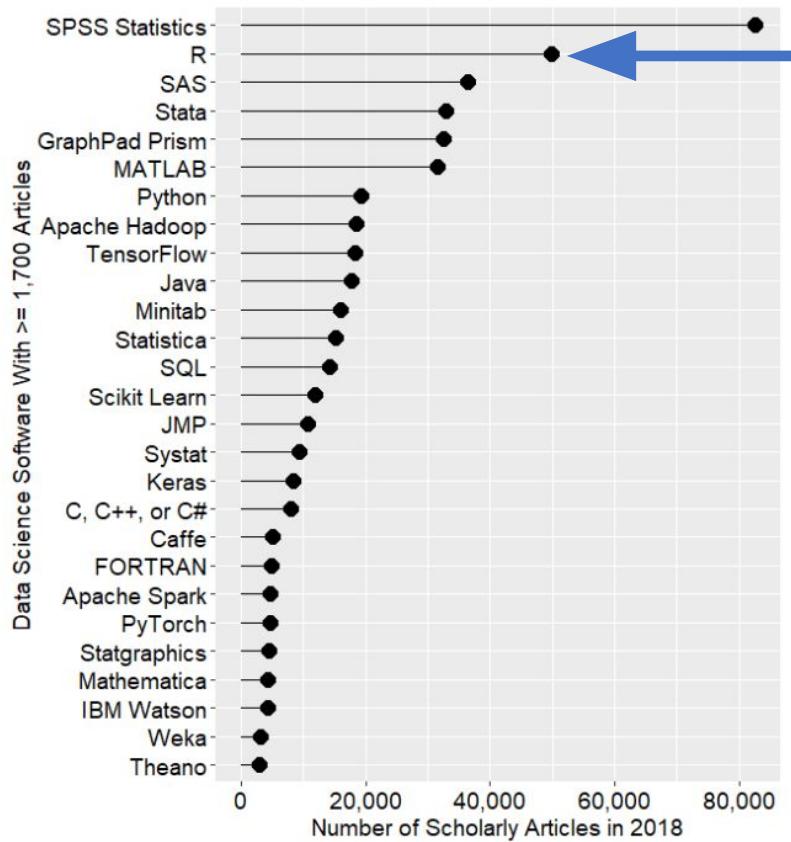


Why learn R?

- R has a large and welcoming community
 - RStudio community
 - Local R User Groups (meetup)
 - R-Ladies
 - R-bloggers
 - The Carpentries
 - Stack Overflow
- Popular in Data Science

R tops other programming languages in

academia

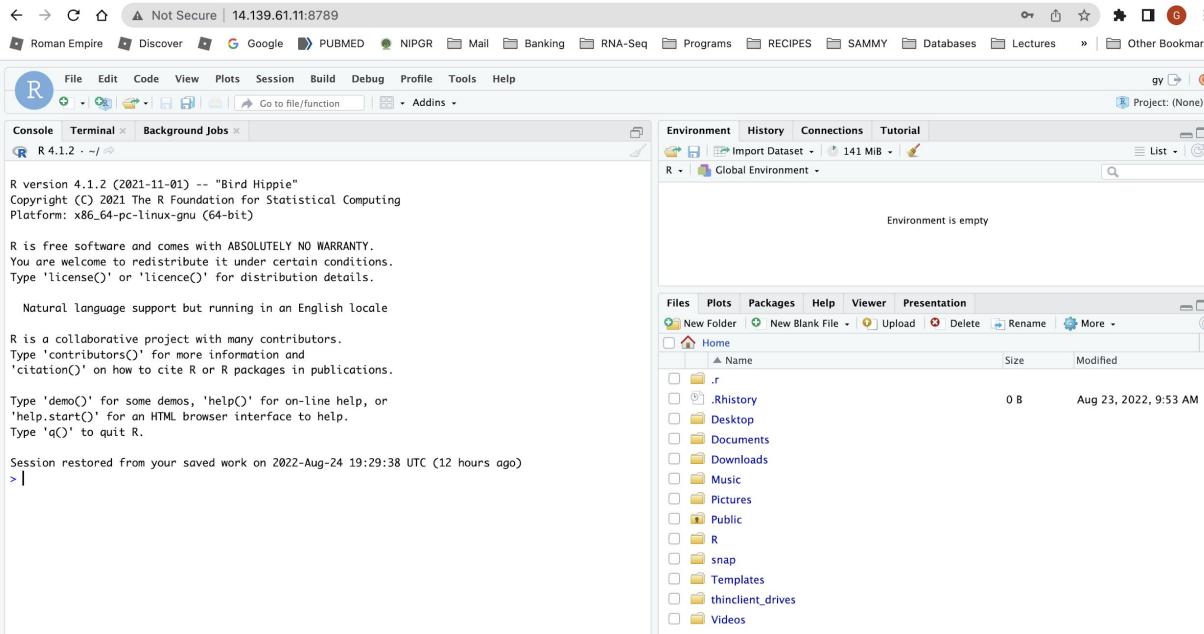




Why learn R?

- Free
- Open-source
- Cross-platform

RStudio



Rstudio IDE : Integrated Development Environment

- An IDE is an application used by software developers that facilitates programming by offering source code editing, building and debugging tools all integrated into one application.
- To function correctly, RStudio needs R and therefore both need to be installed on your computer.
- The RStudio Desktop open-source product is free under the [Afferro General Public License \(AGPL\) v3](#). [Other versions of RStudio](#) are also available.
- We will use RStudio IDE to:
 - write code,
 - navigate files on our computer,
 - inspect the variables we are going to create,
 - visualize the plots we will generate.
- RStudio can also be used for other things (*e.g.*, version control, developing packages, writing Shiny apps) that we will not cover during the workshop.

RStudio

Let us start the action!

Rstudio IDE

The screenshot shows the RStudio IDE interface with three main sections highlighted:

- R script**: The leftmost section contains R code for biomass calculation and estimation. It includes functions like `Biomass calculation per tree` and `plot(kalimantan\$dbh, ...)`.
- R console**: The bottom-left section shows the R console history with commands like `> kai.plot<-merge(kai.plot, Dmed.Hmed.plot, by="Plot")` and `> # calculating the`.
- R environment**: The right section displays the global environment with objects like `hil.trees`, `kal.plot`, and `kalimantan`. A box labeled "R environment" is drawn around this area.
- Graphical output**: The bottom-right section shows a box plot titled "Biomass estimation per plot with different models". The y-axis is "Biomass (Mg?ha⁻¹)" ranging from 100 to 500. The x-axis categories are represented by different symbols. A box labeled "Graphical output" is drawn around this plot.

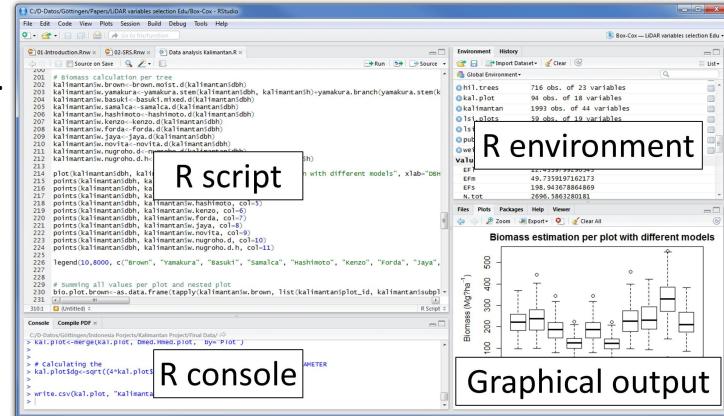
```

C:/D-Datos/Göttingen/Papers/LiDAR variables selection Edu/Box-Cox - RStudio
File Edit Code View Plots Session Build Debug Tools Help
01-introduction.Rnw 02-SRS.Rnw Data analysis Kalimantan.R
Source on Save Run Source
200
201 # Biomass calculation per tree
202 kalimantan$w.brown<-brown.moist.d(kalimantan$dbh)
203 kalimantan$w.yamakura<-yamakura.stem(kalimantan$dbh, kalimantan$h)+yamakura.branch(yamakura.stem(k
204 kalimantan$w.basuki<-basuki.mixed.d(kalimantan$dbh)
205 kalimantan$w.samalca<-samalca.d(kalimantan$dbh)
206 kalimantan$w.hashimoto<-hashimoto.d(kalimantan$dbh)
207 kalimantan$w.kenzo<-kenzo.d(kalimantan$dbh)
208 kalimantan$w.forda<-forda.d(kalimantan$dbh)
209 kalimantan$w.jaya<-jaya.d(kalimantan$dbh)
210 kalimantan$w.novita<-novita.d(kalimantan$dbh)
211 kalimantan$w.nugroho.d<-nugroho.d(kalimantan$dbh)
212 kalimantan$w.nugroho.d.h<-nugroho.d.h(kalimantan$dbh)
213
214 plot(kalimantan$dbh, kalimantan$w.brown, xlab="DBH", ylab="Biomass (Mg?ha-1)", main="Biomass estimation per plot with different models", xlab="DBH", ylab="Biomass (Mg?ha-1)")
215 points(kalimantan$dbh, kalimantan$w.yamakura, col=2)
216 points(kalimantan$dbh, kalimantan$w.basuki, col=3)
217 points(kalimantan$dbh, kalimantan$w.samalca, col=4)
218 points(kalimantan$dbh, kalimantan$w.hashimoto, col=5)
219 points(kalimantan$dbh, kalimantan$w.kenzo, col=6)
220 points(kalimantan$dbh, kalimantan$w.forda, col=7)
221 points(kalimantan$dbh, kalimantan$w.jaya, col=8)
222 points(kalimantan$dbh, kalimantan$w.novita, col=9)
223 points(kalimantan$dbh, kalimantan$w.nugroho.d, col=10)
224 points(kalimantan$dbh, kalimantan$w.nugroho.d.h, col=11)
225
226 legend(10,8000, c("Brown", "Yamakura", "Basuki", "Samalca", "Hashimoto", "Kenzo", "Ford", "Jaya", "Nugroho"))
227
228
229 # Summing all values per plot and nested plot
230 bio.plot.brown<-as.data.frame(tapply(kalimantan$w.brown, list(kalimantan$plot_id, kalimantan$subplot_id), sum))
231
232
310:1 (Untitled) R Script
Console Compile PDF
C:/D-Datos/Göttingen/Indonesia Projects/Kalimantan Project/Final Data/
> kai.plot<-merge(kai.plot, Dmed.Hmed.plot, by="Plot")
>
>
> # calculating the
> kai.plot$dg<-sqrt((4*kai.plot$)
>
>
> write.csv(kai.plot, "Kalimantan.csv")
>

```

Rstudio IDE

- RStudio interface screenshot.
- Clockwise from top left: Source, Environment/History, Files/Plots/Packages/Help/Viewer, Console.
- RStudio is divided into 4 “Panes”
- The **Source/Script** for your scripts and documents (top-left, in the default layout),
- Your **Environment/History** (top-right),
- Your **Files/Plots/Packages/Help/Viewer** (bottom-right)
- R **Console** (bottom-left). The placement of these panes and their content can be customized (see menu, Tools -> Global Options -> Pane Layout).
- **Advantages :**
 - All information you need to write code is available in a single window.
 - Many shortcuts, autocompletion, and highlighting
 - RStudio makes typing easier and less error-prone



INTERACTING WITH R!

1. Write down instructions for the computer to follow (**CODE!**)
2. Tell the computer to follow those instructions. (**EXECUTE/RUN THE CODE!**)

We **write, or code,** instructions in R because it is a common language that both the computer and we can understand!

INTERACTING WITH R!

- There are two main ways of interacting with R:
 - By using the Console
 - By using Script files (plain text files that contain your code).

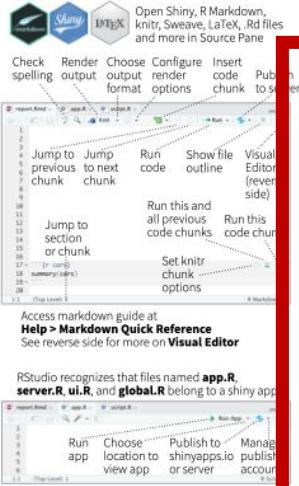
ALL R STUDIO CHEAT SHEETS ARE HERE :

<https://www.rstudio.com/resources/cheatsheets/>

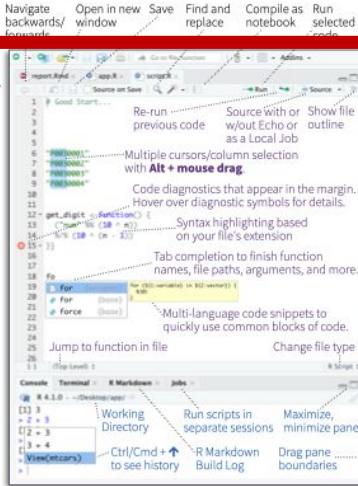
RStudio IDE :: CHEAT SHEET



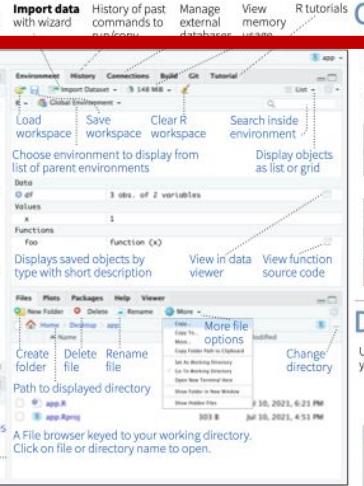
Documents and Apps



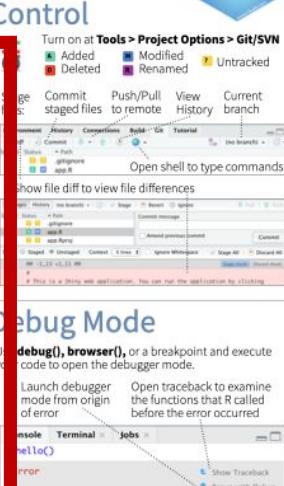
Source Editor



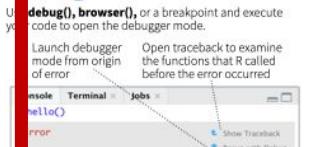
Tab Panes



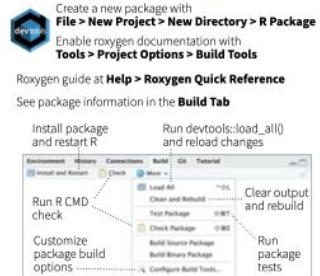
Version Control



Debug Mode



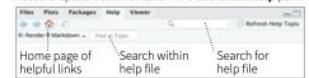
Package Development



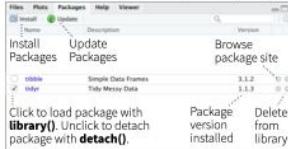
RStudio opens plots in a dedicated **Plots** pane



RStudio opens documentation in a dedicated **Help** pane



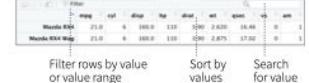
GUI Package manager lists every installed package



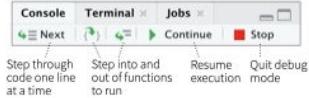
Viewer pane displays HTML content, such as Shiny apps, RMarkdown reports, and interactive visualizations



View(<data>) opens spreadsheet-like view of data set



Run commands in environment where execution has paused
Examine variables in executing environment
Select function in traceback to debug



Source Editor

Navigate backwards/forwards Open in new window Save Find and replace Compile as notebook Run selected code

Re-run previous code
Source with or w/out Echo or as a Local Job
Show file outline
Multiple cursors/column selection with **Alt + mouse drag**.
Code diagnostics that appear in the margin. Hover over diagnostic symbols for details.
Syntax highlighting based on your file's extension
Tab completion to finish function names, file paths, arguments, and more.
Jump to function in file
Multi-language code snippets to quickly use common blocks of code.
Change file type
Working Directory
Run scripts in separate sessions
Maximize, minimize panes
Ctrl/Cmd + ↑ to see history
R Markdown Build Log
Drag pane boundaries

Tab Panes

Import data with wizard History of past commands to run/copy Manage external databases View memory usage R tutorials

Environment History Connections Build Git Tutorial
Import Dataset 148 MB
Global Environment
Load workspace Save workspace Clear R workspace Search inside environment
Choose environment to display from list of parent environments
Display objects as list or grid
Data 3 obs. of 2 variables
x 1
Functions Foo function ()
Displays saved objects by type with short description
View in data viewer View function source code
Files Plots Packages Help Viewer
New Folder Home Desktop App
Create folder Delete file Rename file More
Name
More file options
Copy File To... Move
Copy Folder Path to Clipboard Set As Working Directory
Set As Working Directory Open New Terminal Here
Open New Terminal Here
More Folder in New Window
More Hidden Files
Jul 30, 2023, 6:21 PM Jul 30, 2023, 4:51 PM
Path to displayed directory
Working Directory
Run scripts in separate sessions
Maximize, minimize panes
Drag pane boundaries
A File browser keyed to your working directory. Click on file or directory name to open.

Creating an R Project

- Before starting to write code in RStudio, we need to create an R Project.
- The idea behind an R project is to have a workspace where you can keep all the files and settings associated with the project together.

To create an “R Project”:

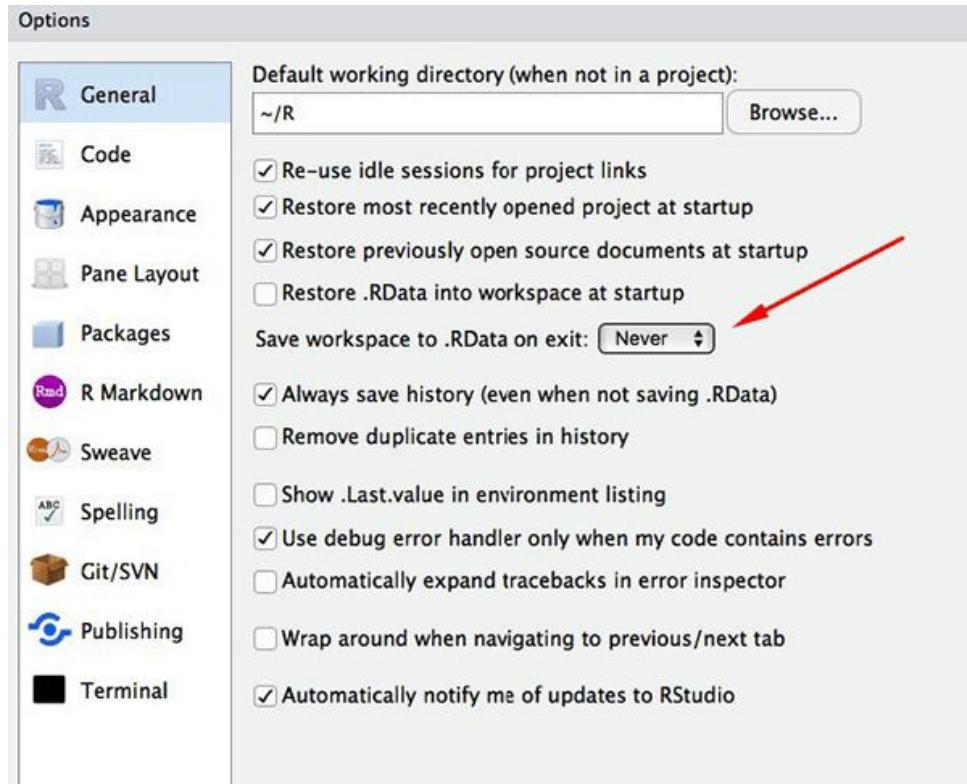
1. Start RStudio.
2. Under the `File` menu, click on `New Project`. Choose `New Directory`, then `New Project`.
3. Enter a name for this new folder (or “directory”), and choose a convenient location for it. This will be your **working directory** for the rest of the day (e.g., `~/my-first-project`).
4. Click on `Create Project`.
5. (Optional) Set Preferences to ‘Never’ save workspace in RStudio.

- Next time you open the R Project it would be easier to resume work.

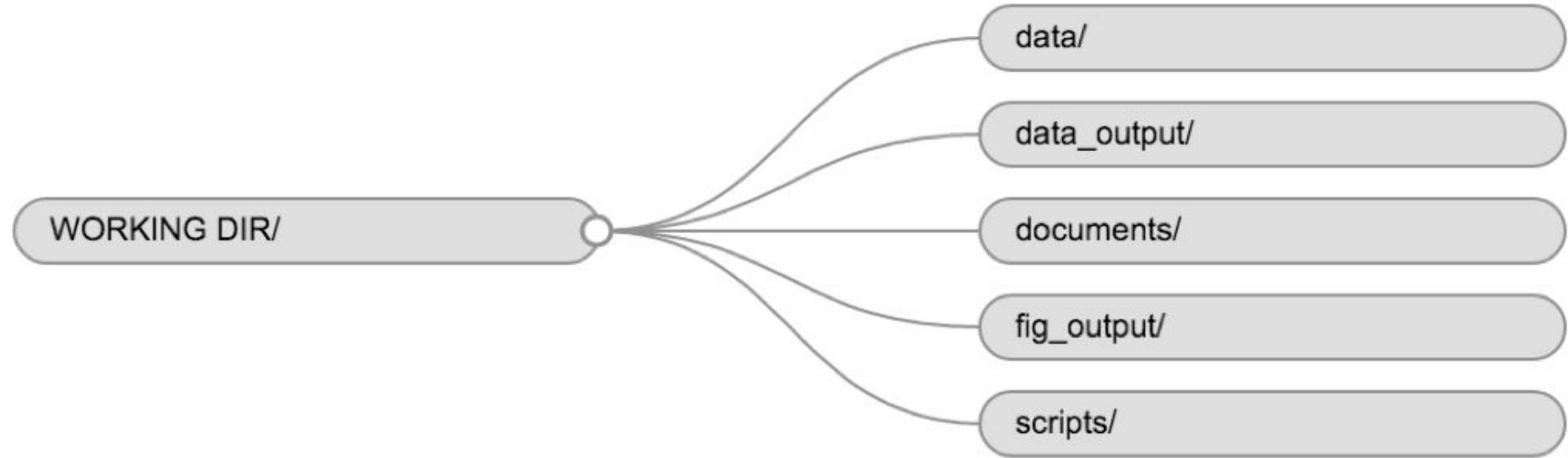


Some Global Options!

- RStudio's default preferences generally work well, but saving a workspace to .RData can be cumbersome, especially if you are working with larger datasets as this would save all the data that is loaded into R into the .RData file.
- To turn that off,
 - go to **Tools** → Global Options
 - Select the 'Never' option for **Save workspace to .RData' on exit.**



Working Directory

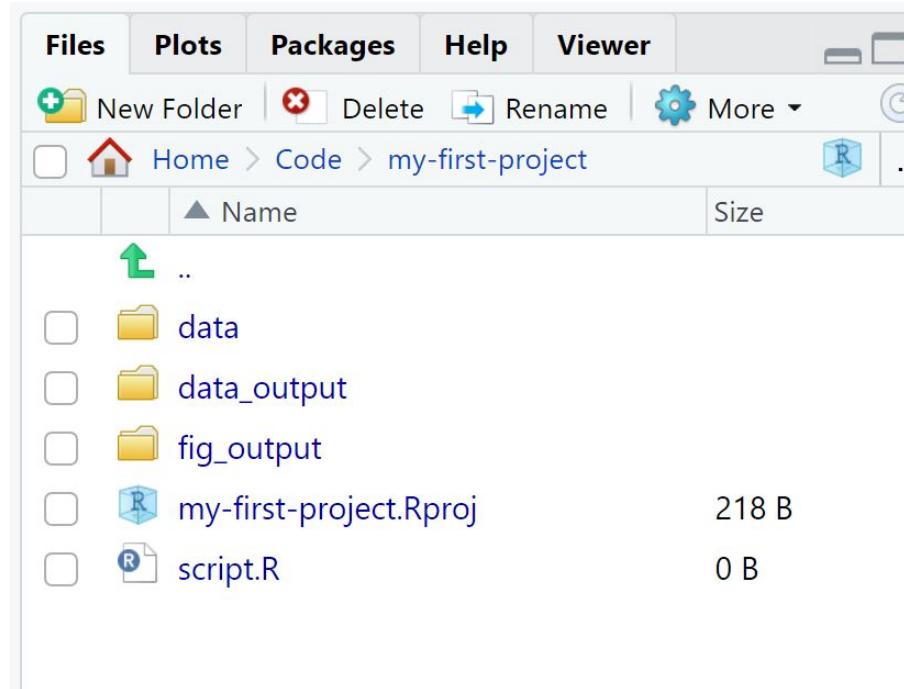


Your Working Directory

- For this workshop, we will need a **data/** folder to store our raw data,
- We will use **data_output/** for when we learn how to export data as CSV files
- We will also create a **fig_output/** folder for the figures that we will save
 - Under the `Files` tab on the right of the screen, click on `New Folder` and create a folder named `data` within your newly created working directory (e.g., `~/my-first-project/data`). Repeat these operations to create a `data_output/` and a `fig_output` folders.
- We are going to keep the **script** in the root of our **working directory** because we are only going to use one file and it will make things easier.



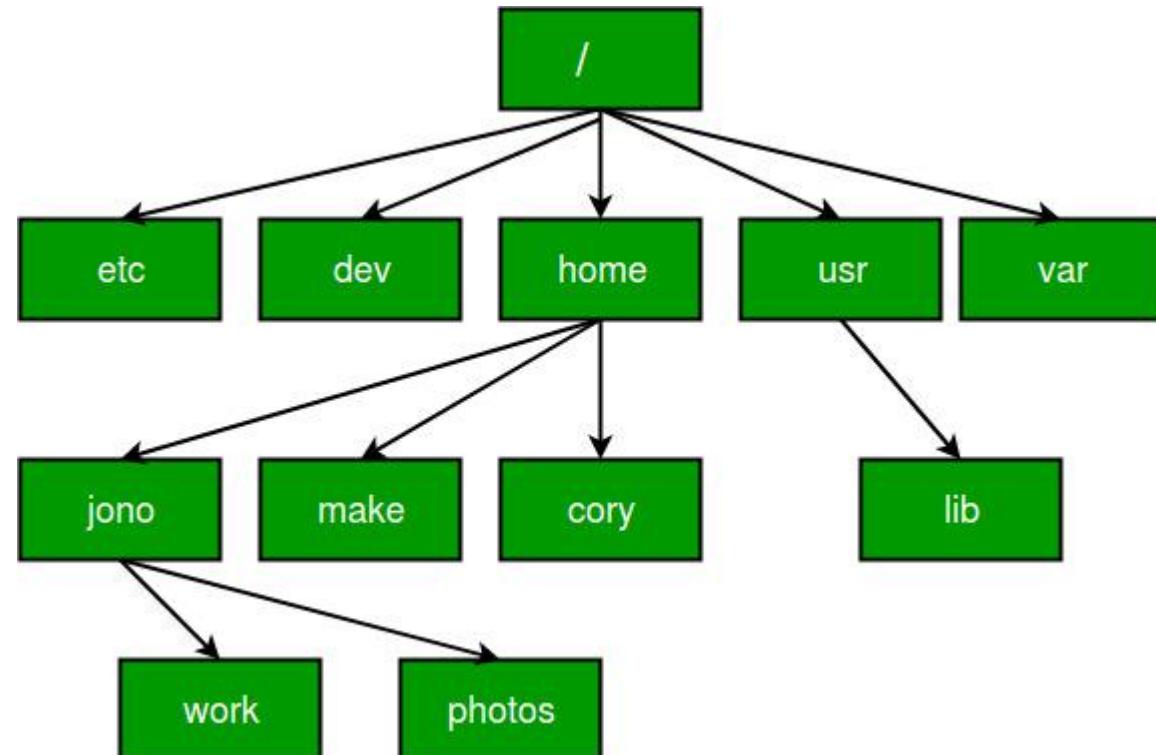
Your Working Directory



Computing Environments

Concept of Working Directory

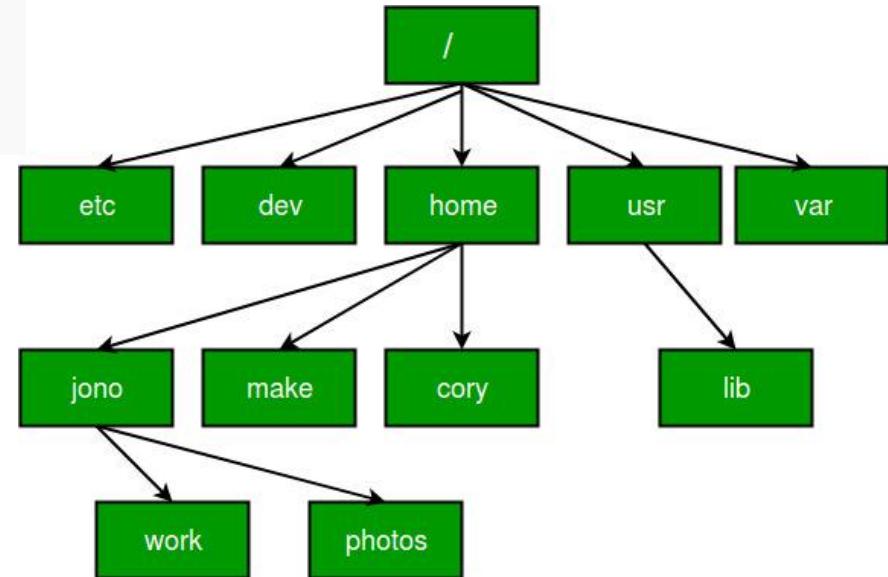
- The working directory is an important concept to understand.
- It is the place from where R will be looking for and saving the files.
- When you write code for your project, it should refer to files in relation to the root of your working directory and only need files within this structure.



WORKING DIRECTORY

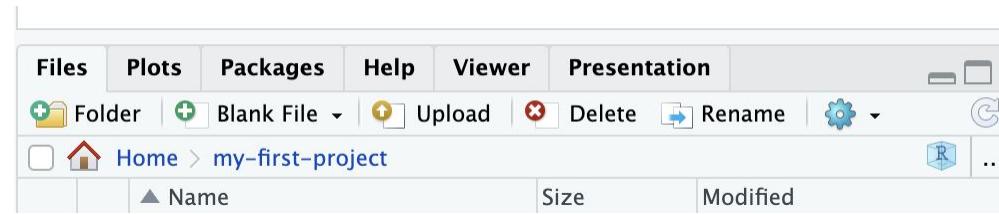
An **absolute path** is defined as specifying the location of a file or directory from the root directory (/). In other words, we can say that an absolute path is a complete path from start of actual file system from / directory.

- Working this way makes it a lot easier to move your project around on your computer and share it with others without worrying about whether or not the underlying scripts will still work!



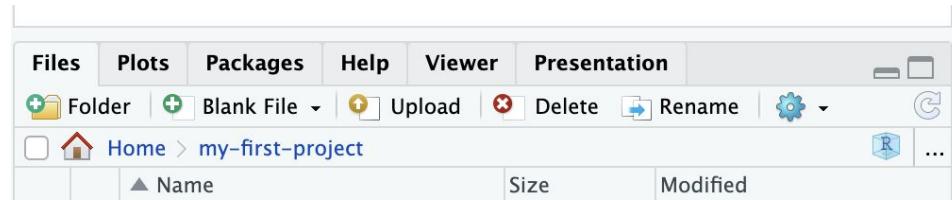
WHAT IS MY WORKING DIRECTORY?!

- Using RStudio projects makes it easy to organise your files in the project and ensures that your working directory is set properly.
- RStudio shows your current working directory at the top of your window:



WHAT IS MY WORKING DIRECTORY?!

- Another way to check your WD is
 - by typing `getwd()` in the console pane.
 - **TRY IT NOW!**
- If for some reason your working directory is not what it should be, you can change it in the RStudio interface by navigating in the file browser where your working directory should be, and clicking on the blue gear icon “More”, and select “Set As Working Directory”.
- Alternatively you can use `setwd("/path/to/working/directory")` to reset your working directory (not recommended). However, your scripts should not include this line because it will fail on someone else’s computer.

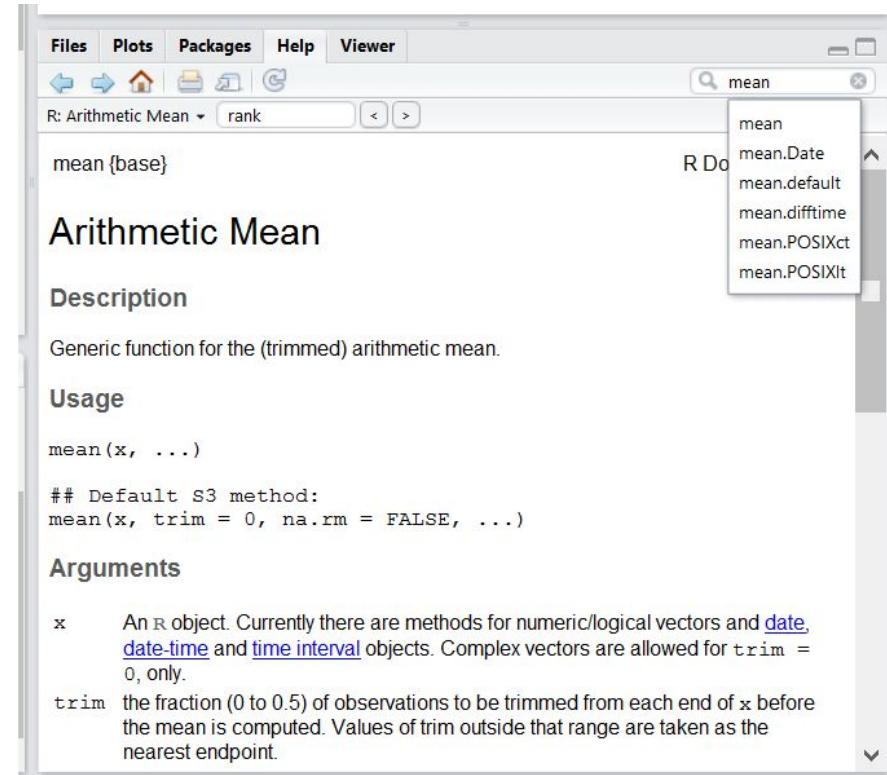


SEEKING HELP

RSTUDIO HELP INTERFACE

Use this to search for
more information on R
functions

TRY IT
NOW!



SEEKING HELP

I know the name of the function I want to use, but I'm not sure how to use it

If you need help with a specific function, let's say `barplot()`, you can type:

```
?barplot
```

**TRY IT
NOW!**



UNIVERSITY OF
CAMBRIDGE

SEEKING HELP

I know the name of the function but I don't know how many / What inputs it requires to run?

If you just need to remind yourself of the names of the arguments, you can use:

```
args(lm)
```

TRY IT NOW!

SEEKING HELP

I want to use a function that does X, there must be a function for it but I don't know which one...

If you are looking for a function to do a particular task, you can use the `help.search()` function, which is called by the double question mark `??`. However, this only looks through the installed packages for help pages with a match to your search request

```
??kruskal
```



SEEKING HELP

I am stuck... I get an error message that I don't understand

Start by googling the error message. However, this doesn't always work very well because often, package developers rely on the error catching provided by R. You end up with general error messages that might not be very helpful to diagnose a problem (e.g. "subscript out of bounds"). If the message is very generic, you might also include the name of the function or package you're using in your query.

Asking for help

The key to receiving help from someone is for them to rapidly grasp your problem. You should make it as easy as possible to pinpoint where the issue might be.

PROGRAMMING IN R!

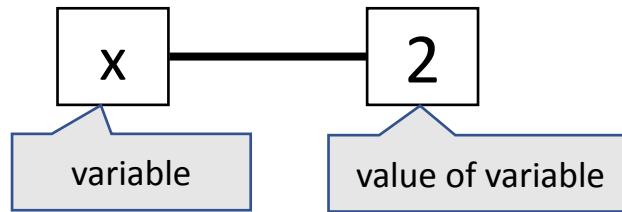
CREATING OBJECTS IN R!

CREATING OBJECTS IN R!

- **Type Math in the Console?**
 - Fun- but not useful!
- **Useful Stuff gets done when you can**
 - Assign Values to Objects!

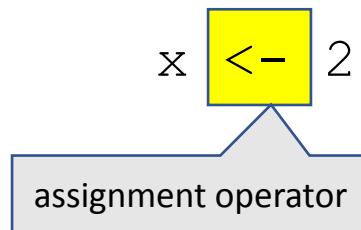
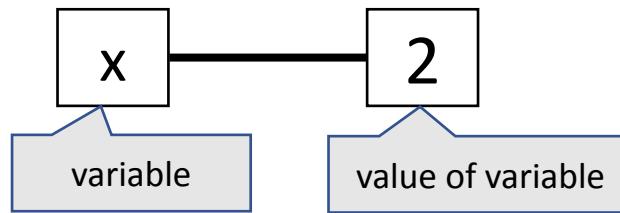
Variables

A variable is a letter or word that stores a value in it.



Variables

A variable is a letter or word that stores a value in it.



Variable names

- Case-sensitive
- Do not start with numbers
- Not too long
- Some names cannot be used as they are reserved
e.g. `if <- 66` 
- Do not use function names
- Use nouns if using a word to represent a variable
e.g. `gene_id <- 12345` 
- Be consistent - style guide



A SMALL EXERCISE!

Challenge

What are the values after each statement in the following?

```
mass <- 47.5          # mass?  
age  <- 122           # age?  
mass <- mass * 2.0    # mass?  
age  <- age - 20      # age?  
mass_index <- mass/age # mass_index?
```



A SMALL EXERCISE!

▼ Answer

```
mass <- 47.5          # mass is 47.5
age  <- 122           # age is 122
mass <- mass * 2.0    # mass is 95
age  <- age - 20      # age is 102
mass_index <- mass/age # mass_index is 0.931
```



FUNCTIONS & THEIR ARGUMENTS

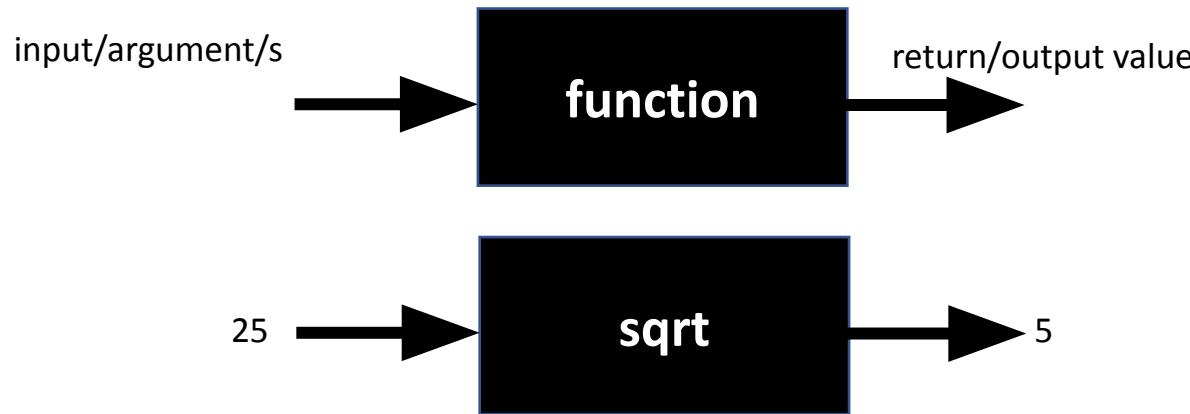
Calling functions

Functions execute a defined set of commands – automate a process



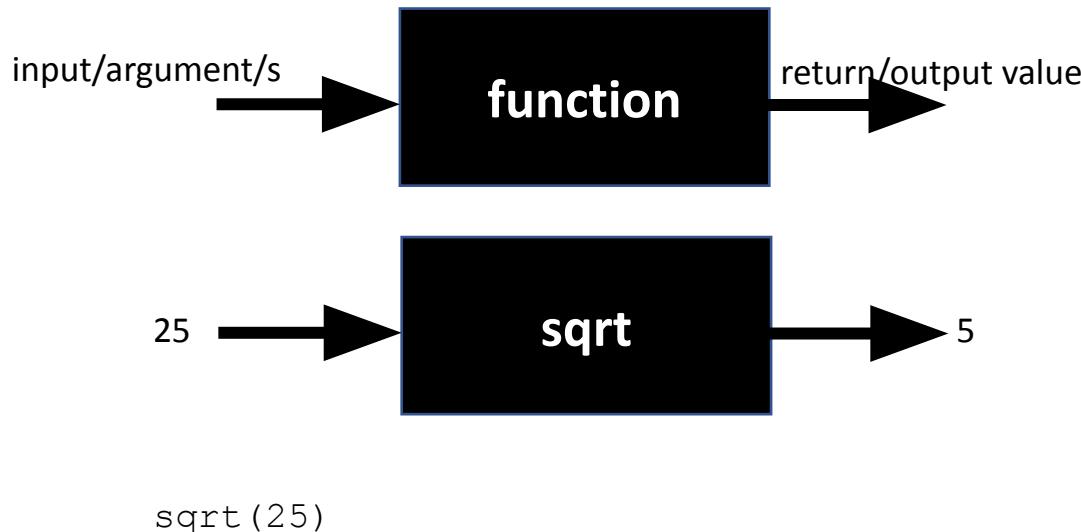
Calling functions

Functions execute a defined set of commands – automate a process



Calling functions

Functions execute a defined set of commands – automate a process



VECTORS & DATA TYPES

Data types

- **logical:** TRUE FALSE
- **integer:** whole numbers. *e.g.*, 40
- **double:** numbers with decimal points. *e.g.*, 2 . 666
- **character:** words or **strings**. *e.g.*, "Hello"

Data structures

- **vector**: list of items of the same data type. *e.g.*, 4, 6, 9, 12
- **factor**: categorical data (has to be a character vector). *e.g.*, Male, Female
- **data.frame**: contains tabular data – normally data is loaded into data.frame when reading in a file

Vector

- List of data types (must be same type)
- One-dimensional

4	12	7	9
---	----	---	---

```
c(4,12,7,9)
```



SUBSETTING VECTORS

Vector indices

1	2	3	4
4	12	7	9

CHALLENGE

Challenge

1. Create a vector `height_mm` that contains the following heights of animals 100,150,99,87 in mm respectively.

▶ Answer

2. What is the total height of these animals in mm? (*Hint:* use function `sum`)

▶ Answer

3. Another animal was measured and its height is 220mm. Add this to the beginning of the vector `height_mm` and save the total height of the new set of animals in the `total_height` variable. What is the new `total_height`?

▶ Answer



FUN QUIZ!

- What is temperature?

temperature <- 26.789

round (x = temperature, digits = 1)

- A. Variable
- B. Function
- C. Place holder

(Poll Question – Choose the right option!)

FUN QUIZ !

- What is the value of temperature after executing both lines of code?

```
temperature <- 26.789
```

```
round (x = temperature, digits = 1)
```

- A. 27
- B. 26.7
- C. 26.789
- D. 26.8

(Poll Question – Choose the right option!)



UNIVERSITY OF
CAMBRIDGE

CONDITIONAL SUBSETTING

Relational Operators

Used to compute a condition/comparisons:

`==` equal

`>` greater than

`<` less than

`>=` greater than or equal to

`<=` less than or equal to

`!=` not equal to



Logical operators

& and

| or

! not



AND &

- TRUE & TRUE results in TRUE
- TRUE & FALSE results in FALSE
- FALSE & TRUE results in FALSE
- FALSE & FALSE results in FALSE



OR |

- TRUE | TRUE results in TRUE
- TRUE | FALSE results in TRUE
- FALSE | TRUE results in TRUE
- FALSE | FALSE results in FALSE



NOT !

- `!TRUE` results in `FALSE`
- `!FALSE` results in `TRUE`



HANDLING MISSING DATA

AN EXERCISE

Challenge

1. Using this vector of heights in inches, create a new vector, `heights_no_na`, with the NAs removed.

```
heights <- c(63, 69, 60, 65, NA, 68, 61, 70, 61, 59, 64, 69, 63, 63, NA,
```

2. Use the function `median()` to calculate the median of the `heights` vector.
3. How many people in the set are taller than 67 inches.

▶ Answer

COPY FROM CHAT BOX!: `heights <- c(63, 69, 60, 65, NA, 68, 61, 70, 61, 59, 64, 69, 63, 63, NA, 72, 65, 64, 70, 63, 65)`



UNIVERSITY OF
CAMBRIDGE

AN EXERCISE

▼ Answer

```
heights <- c(63, 69, 60, 65, NA, 68, 61, 70, 61, 59, 64, 69, 63, 63)

# 1.
heights_no_na <- heights[!is.na(heights)]
# or
heights_no_na <- na.omit(heights)
# or
heights_no_na <- heights[complete.cases(heights)]

# 2.
median(heights, na.rm = TRUE)

# 3.
heights_above_67 <- heights_no_na[heights_no_na > 67]
length(heights_above_67)
```



What we have learned so far

- How to create a Project in RStudio
- How to code and execute R code in RStudio
- Create variables *e.g.*, `weight <- 24`
- Data types: integer (11), double (11.01), character ("Hello"), logical (TRUE/FALSE)
- Calling functions *e.g.*, `sqrt(25)`
- Create vectors *e.g.*, `weight_mm <- c(22, 24, 10, 34)`
- Index vectors *e.g.*, `weight_mm[3]`
- Relational and Logical operators (`& | ! == < > <= >= !=`)
- Missing data `NA`



UP NEXT

STARTING WITH REAL DATA!

