# UNIVERSITY OF CAMBRIDGE

Bioinformatics Training

# Introduction to R
## for Biologists

Alison Horst Artwork!

**Gita Yadav**

# Introductions

Tutors

Course participants introductions

UNIVERSITY OF
CAMBRIDGE

# This course

- Aimed for **beginners** who never used R before.

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

UNIVERSITY OF CAMBRIDGE

# Course Materials and Links

Online RStudio access link

TO COVER THIS MORNING

- RSTUDIO IDE
- CREATING PROJECTS & SCRIPTS
- VARIABLES
- FUNCTIONS
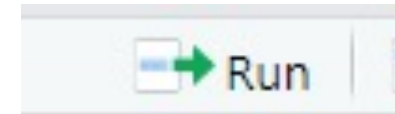- DATA TYPES / DATA STRUCTURES
- VECTORS
- INDEXING

UNIVERSITY OF CAMBRIDGE

# Why learn R?

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

**VS**

**Rscript <name of file>**

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.

UNIVERSITY OF
CAMBRIDGE

# Why learn R?

- R is **interdisciplinary**
  - Over 10,000 packages that cover different fields
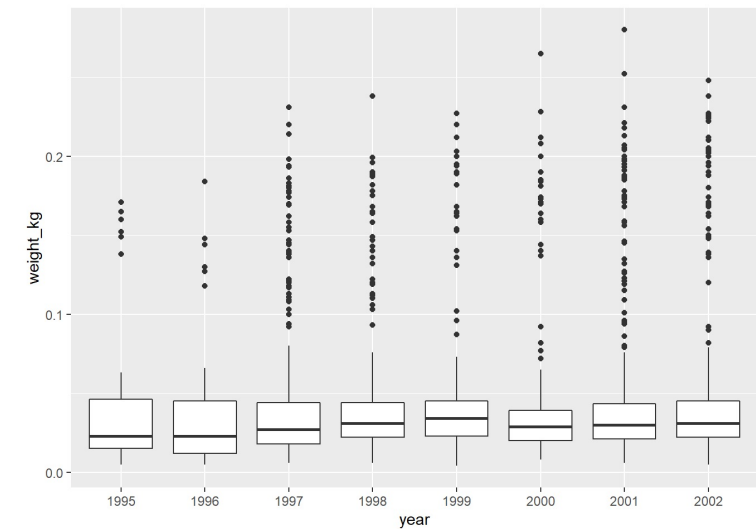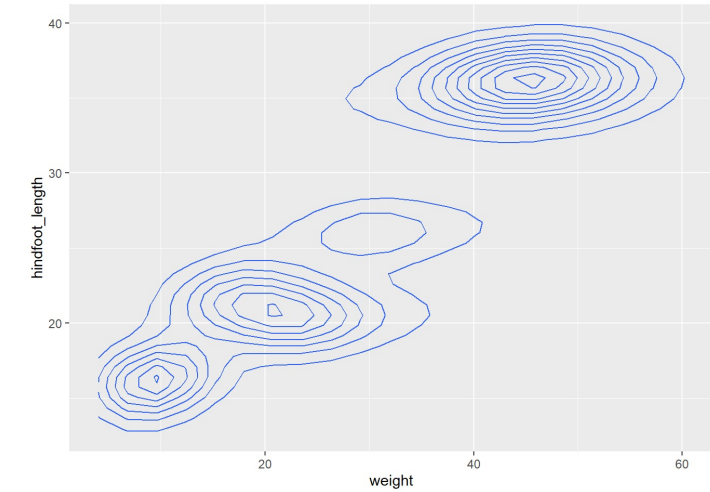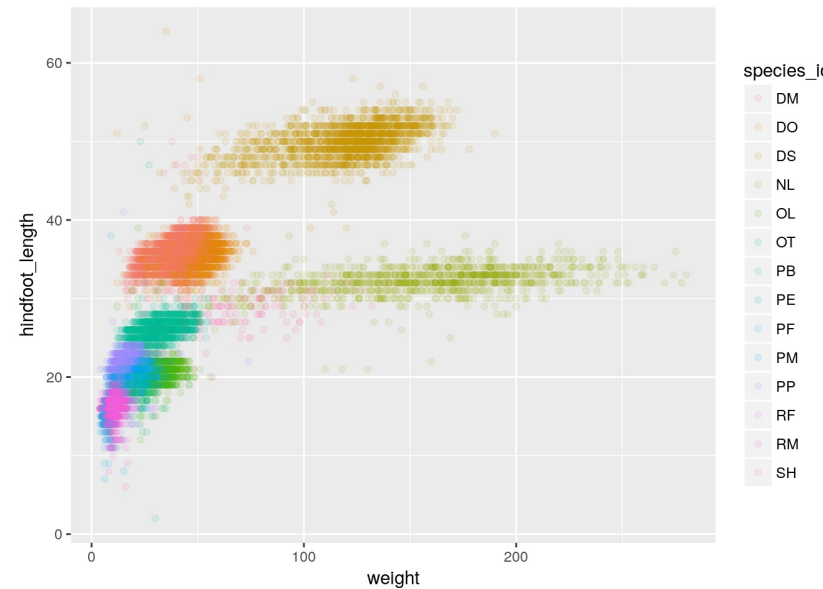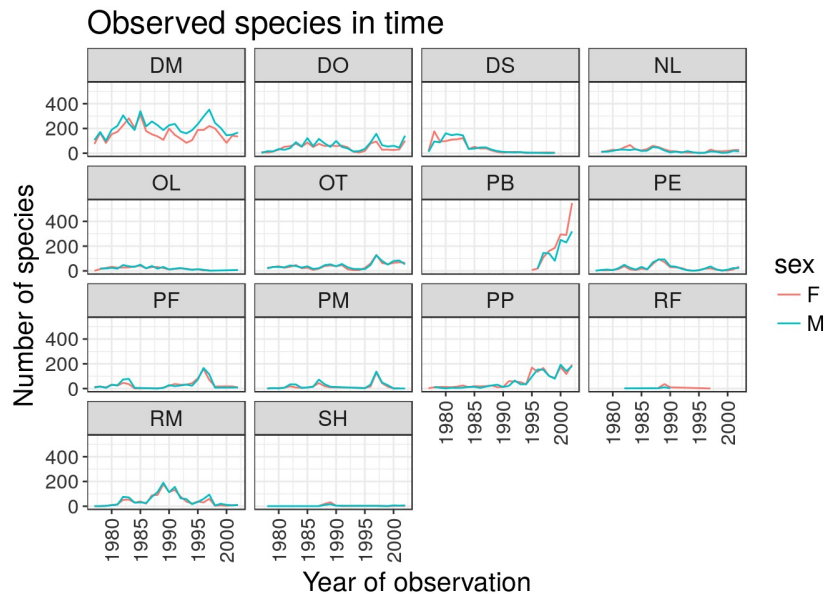  - Bioconductor repository for bioinformatics packages

UNIVERSITY OF CAMBRIDGE

# Why learn R?

- R works on data of different sizes

UNIVERSITY OF CAMBRIDGE
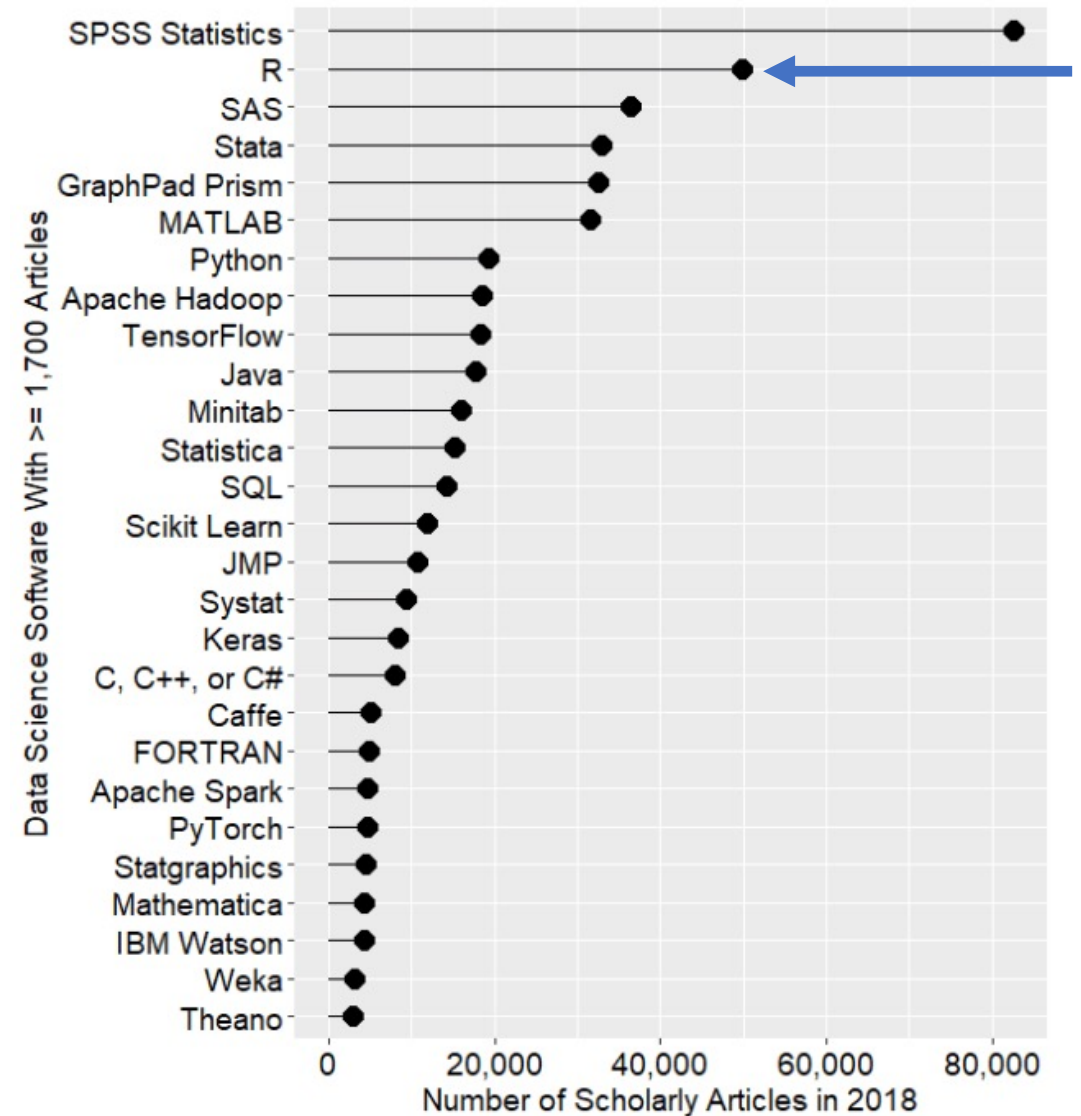
# Why learn R?

- R produces high-quality graphics.

UNIVERSITY OF CAMBRIDGE

# 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

UNIVERSITY OF CAMBRIDGE

# Why learn R?

- Free

- Open-source

- Cross-platform

UNIVERSITY OF
CAMBRIDGE

# RStudio

# Let us start the action!

UNIVERSITY OF CAMBRIDGE

# Rstudio IDE

UNIVERSITY OF CAMBRIDGE

# Working Directory

# Variables

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

UNIVERSITY OF
CAMBRIDGE

# 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 <u>reserved</u>

  *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 - <u>style guide</u>

Bioinformatics Training

UNIVERSITY OF
CAMBRIDGE

# Calling functions

Functions execute a defined set of commands – automate a process



input/argument/s → **function** → return/output value

UNIVERSITY OF CAMBRIDGE

# Calling functions

Functions execute a defined set of commands – automate a process

input/argument/s → **function** → return/output value

25 → **sqrt** → 5

UNIVERSITY OF CAMBRIDGE

# Calling functions

Functions execute a defined set of commands – automate a process



```
sqrt(25)
```

UNIVERSITY OF CAMBRIDGE

# 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!)**

UNIVERSITY OF CAMBRIDGE

# 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!)**

Bioinformatics Training

UNIVERSITY OF CAMBRIDGE

# 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"`

UNIVERSITY OF
CAMBRIDGE

# 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

UNIVERSITY OF CAMBRIDGE

# Vector

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

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

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

UNIVERSITY OF
CAMBRIDGE

# Vector indices

UNIVERSITY OF CAMBRIDGE

# 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

UNIVERSITY OF
CAMBRIDGE

# NOT !

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

# 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`

Bioinformatics Training

UNIVERSITY OF CAMBRIDGE