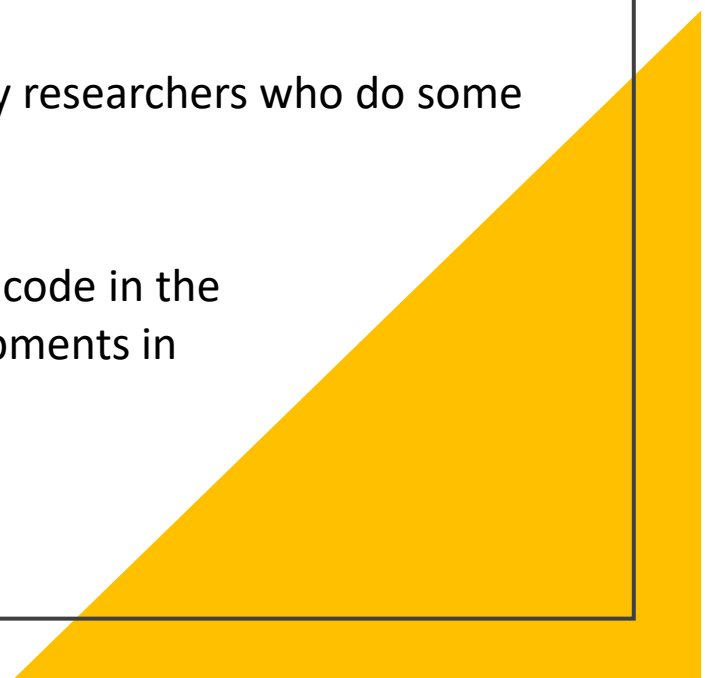# 15$^{th}$ Bioinformatics in Health Sciences course

# Introduction to R

**Ana Gonçalves**
*ICVS/3B's, University of Minho*

**Nuno S. Osório**
*ICVS/3B's, University of Minho*

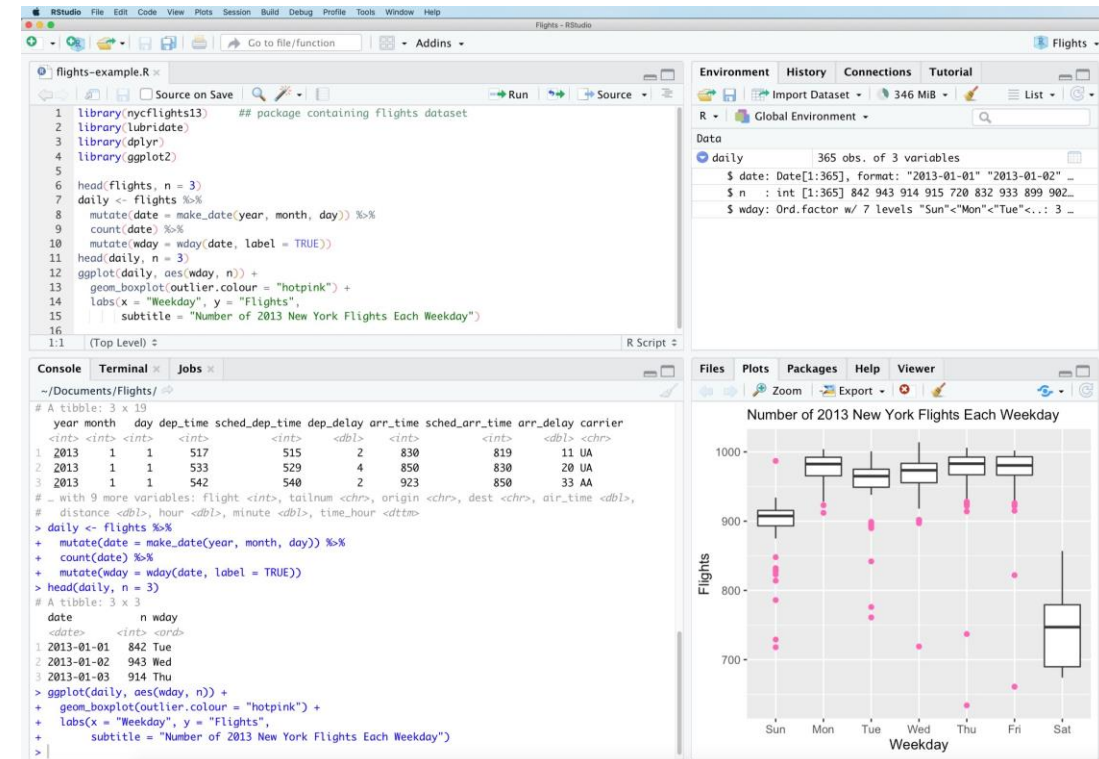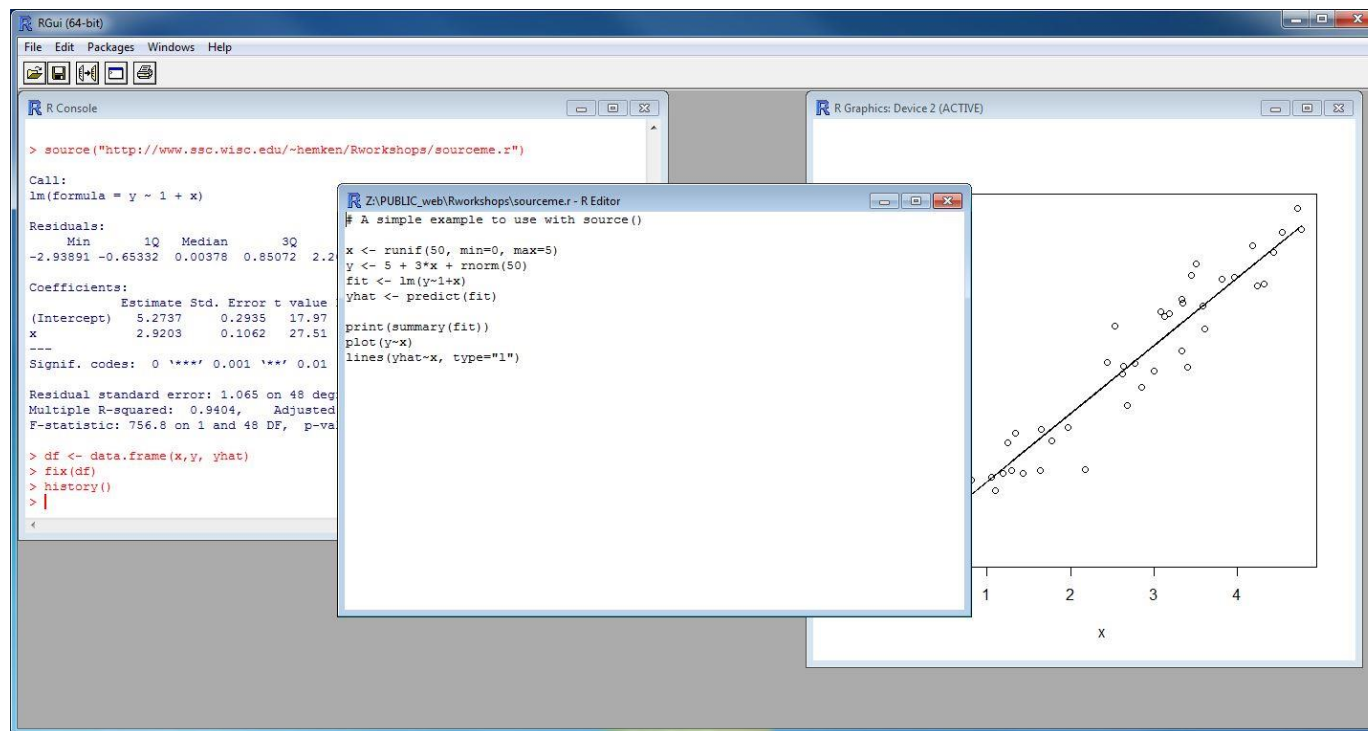February 7$^{th}$ 2024

# WHY LEARN R?

- R is available for Windows, OS X, and Linux/Unix platforms.

- R is free and open-source software, with no license restrictions.

- R plays a key role in a wide variety of research and data analysis projects because it makes many modern statistical methods, both simple and advanced, readily available and easy to use.

- R provides a wonderfully flexible programming environment favored by the many researchers who do some form of data analysis as part of their work.

- Perhaps the most appealing feature of R is that any programmers can contribute code in the form of *packages* (or *libraries*), so the rest of the world has fast access to developments in statistics and data science.

**VS**

R refers to a programming language as well as the software that runs R code.
RStudio is a software interface that can make it easier to write R scripts and interact with the R software.

# Time to do some coding!

**Let's go open both R and Rstudio to further understand the differences!**

Sciences course15$^{th}$ Bioinformatics in Health
Introduction to R

February 7$^{th}$ 2024

ICVS

University of Minho
School of Medicine

ICVS/3B's
Associate
Laboratory
University of Minho

# R: Packages

Packages are the fundamental units of reproducible R code. They include reusable R functions, the documentation that describes how to use them, and sample data.

- Once a package is installed, we don't need to install it again in our R studio.

## BUT...

- **We always need to call a library in order to run specific functions.**

```
63  library(BiocGenerics)
64  library(BiocManager)
65  library(readxl)
66  library(ggplot2)
67  library(gprofiler2)
68  library(clustifyr)
69  library(org.Hs.eg.db)
70  library(org.Mm.eg.db)
71  library(clusterProfiler)
72  library(AnnotationDbi)
73  library(dplyr)
74  library(preprocessCore)
75  library(DEGreport)
76  library(enrichplot)
77  library(GOSemSim)
78  library(vidger)
79  library(R.utils)
80  library(Biobase)
```

```
29
30  ```{r}
31
32  if (!require("BiocManager", quietly = TRUE))
33      install.packages("BiocManager")
34
35
36  BiocManager::install("clusterProfiler")
37  BiocManager::install("enrichplot")
38
39  BiocManager::install("DESeq2")
40  BiocManager::install("org.Hs.eg.db")
41  BiocManager::install("org.Mm.eg.db")
42  BiocManager::install("AnnotationDbi")
43
44  install.packages("gprofiler2")
45
46  BiocManager::install("clustifyr")
47  BiocManager::install("Biobase")
48  BiocManager::install("GOSemSim")
49  BiocManager::install("vidger")
50  BiocManager::install("DEGreport")
51  BiocManager::install("preprocessCore")
52
53  install.packages("R.utils")
54  install.packages("readxl")
55  ```
```

# R: Packages - BiocManager



https://www.bioconductor.org/

# Open source software for Bioinformatics

The Bioconductor project aims to develop and share open source software for precise and repeatable analysis of biological data. We foster an inclusive and collaborative community of developers and data scientists.

https://www.bioconductor.org/packages/release/bioc/

Home > BiocViews

**Bioconductor version 3.18 (Release)**

Find biocViews:

Software (2266)

# R: files

**To open a .txt file:**

```
# Install utils package
install.packages(" utils ")

# Access utils package
library(utils)
```

Now, depending if the txt file has decimals, a header, and other factors, three functions can be selected:

```
# Read data from .xlsx sheet called "Year1" as data frame and assign to object

1. Data <- read.delim(file = "my_file.txt", header = TRUE, sep = "t", dec = ".")

2. Data <- read.delim2(file = "my_file.txt", header = TRUE, sep = "t", dec = ",")

3. Data <- read.table(file = "my_file.txt", header = TRUE)
```

# R: files

## To open a CSV file:

```
# Install readr package
install.packages("readr")

# Access readr package
library(readr)

# Read .csv file into R as data frame
Data <- read_csv("transcriptomics.csv")

# Read .csv data file into R from a website
data <- read_csv("https://raw.githubusercontent.com/tutorial/transcriptomics.csv ")
```

**We will use this method in our script**

# R: files

**To open a .xlsx file:**

```
# Install readxl package
install.packages("readxl")

# Access readxl package
library(readxl)

# Read data from .xlsx sheet called "transcriptomics" as data frame and assign to object
Data <- read_excel("data.xlsx", sheet="transcriptomics")
```

# Time to do some coding!

**R** Studio®
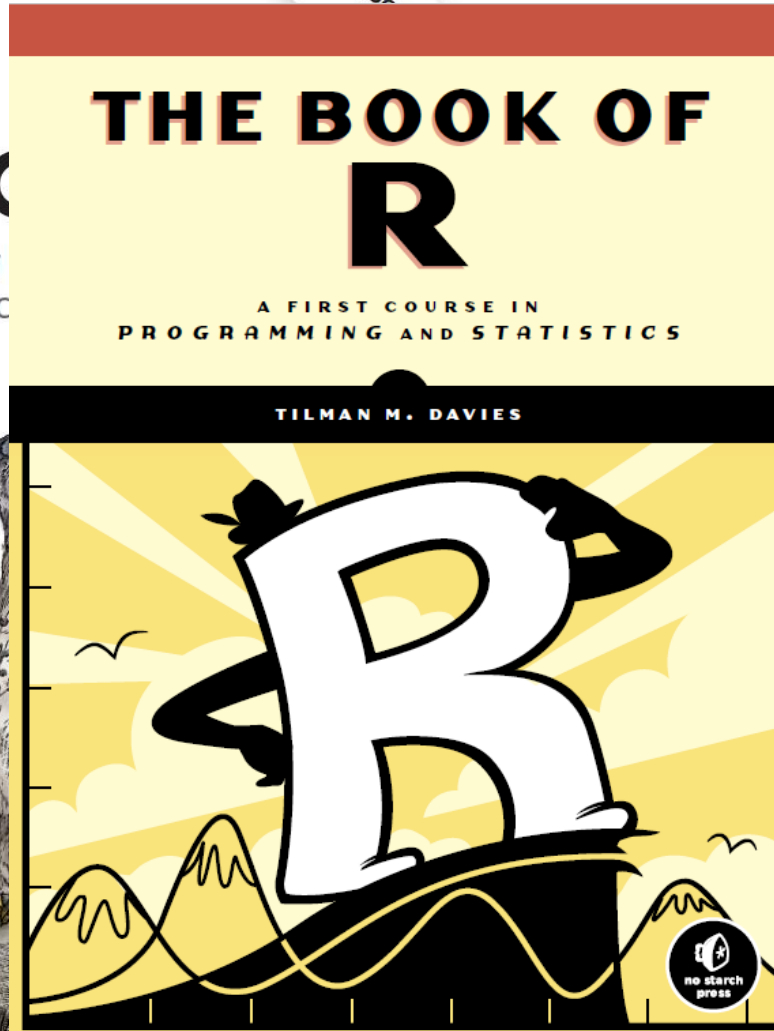
**Sciences course15**[th] **Bioinformatics in Health**
**Introduction to R**

February 7[th] 2024

ICVS

ICVS/3B's

# R: Books recommendations

O'REILLY®

R Coo

Proven Recipes for
Statistics & Graphic

THE BOOK OF
R

A FIRST COURSE IN
PROGRAMMING AND STATISTICS

TILMAN M. DAVIES

no starch press

```
fib <- c(0, 1, 1, 2, 3, 5, 8, 13, 21, 34)
cat("The first few Fibonacci numbers are:", fib, "...\n")
#> The first few Fibonacci numbers are: 0 1 1 2 3 5 8 13 21 34 ...
```

Using cat gives you more control over your output, which makes it especially useful in R scripts that generate output consumed by others. A serious limitation, however, is that it cannot print compound data structures such as matrices and lists. Trying to cat them only produces another mind-numbing message:

```
cat(list("a", "b", "c"))
#> Error in cat(list("a", "b", "c")): argument 1 (type 'list') cannot
#>   be handled by 'cat'
```

### See Also

See Recipe 4.2 for controlling output format.

## 2.2 Setting Variables

### Problem

You want to save a value in a variable.

### Solution

Use the assignment operator (<-). There is no need to declare your variable first:
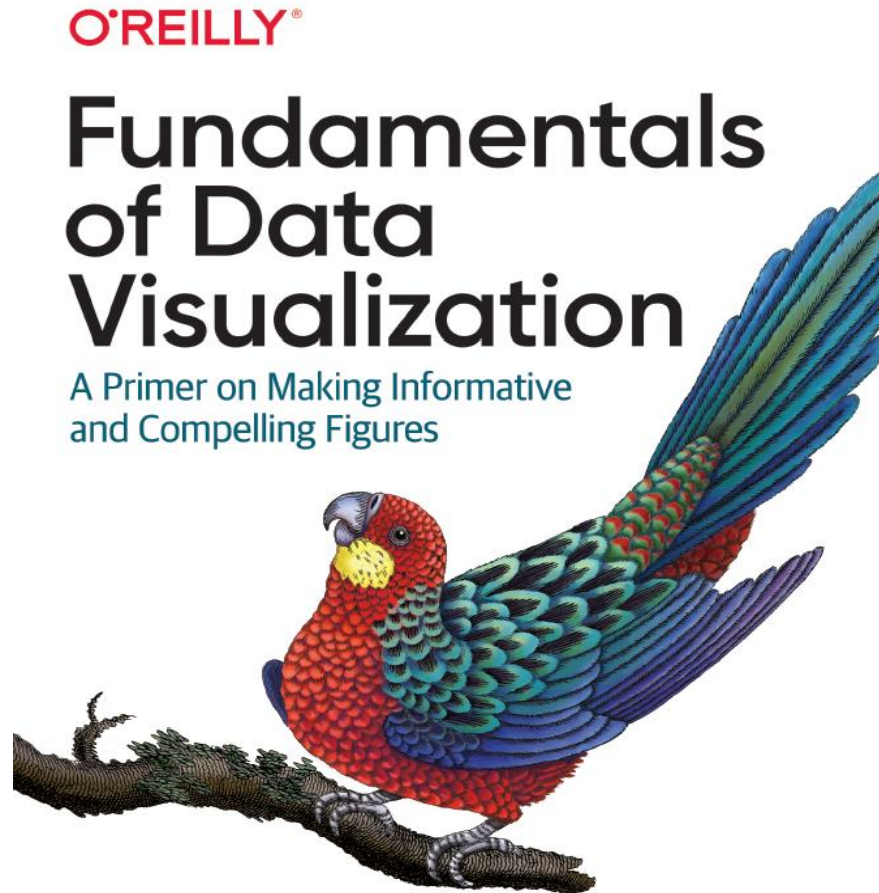
```
x <- 3
```

### Discussion

Using R in "calculator mode" gets old pretty fast. Soon you will want to define variables and save values in them. This reduces typing, saves time, and clarifies your work.

There is no need to declare or explicitly create variables in R. Just assign a value to the name and R will create the variable:
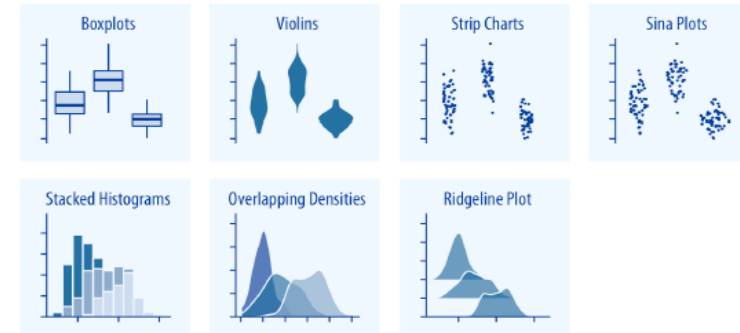
```
x <- 3
y <- 4
z <- sqrt(x^2 + y^2)
print(z)
#> [1] 5
```

Notice that the assignment operator is formed from a less-than character (<) and a hyphen (-) with no space between them.
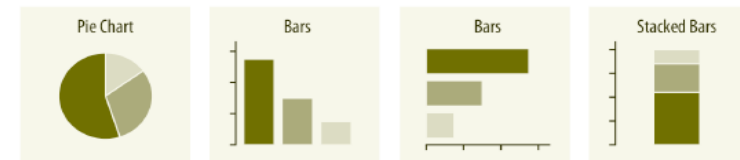
# R: Books recommendations

# R: ggplot2 Package

https://cran.r-project.org/web/packages/ggplot2/index.html



https://r-charts.com/ggplot2/

https://r-graph-gallery.com/ggplot2-package.html

# Time to do some coding!

Sciences course15$^{th}$ Bioinformatics in Health
Introduction to R

February 7$^{th}$ 2024

# 15$^{th}$ Bioinformatics in Health Sciences course

## Using R for Gene Expression Analysis

**Ana Gonçalves**
*ICVS/3B's, University of Minho*

**Nuno S. Osório**
*ICVS/3B's, University of Minho*

February 7$^{th}$ 2024

**ICVS**

University of Minho
School of Medicine

ICVS/3B's
Associate Laboratory
University of Minho

# R: Package for DE genes analysis

**DESeq2:** **A Popular RNA-Seq Analysis Package**

**Widely used and well-documented:** DESeq2 is a widely used tool in the bioinformatics community, with a large user base and a well-documented user manual.

https://bioconductor.org/packages/devel/bioc/vignettes/DESeq2/inst/doc/DESeq2.html

**Workflow**

**Count matrix** → DESeq Object

An object used to store the input values, intermediate calculations and results of the analysis of differential expression.

The reads from the RNA sequencing experiment

↓

Filtering

**DESeq() function**

GLM approach
Love, M.I., Huber, W., Anders, S. (2014)

Normalization: estimateSizeFactors()
Dispersion: estimateDispersions()
Wald statistics: nbinomWaldTest()

Removal of genes that have null reads across all samples

Acquire the differential expressed genes

↓

**DE results**

# R: Information necessary for the DESeq Object

**Count Matrix**

|  | Sample 1 | Sample 2 |  | Sample M |
|---|---|---|---|---|
|  | GSM6160812 | GSM6160813 |  | GSM6160833 |
| Gene 1 | 88 | 10 | ... | 102 |
| Gene 2 | 55 | 33 | ... | 34 |
| ... | ... | ... | ... | ... |
| Gene N | 32 | 0 | ... | 22 |

**Study information** → **Study design**

|  | Condition 1 | Condition 2 | Condition 3 | Condition 4 |
|---|---|---|---|---|
|  | Treatment | Time Point | Sex | Batch |
| GSM6160812 | Placebo | Pre-clinical | Female | 1 |
| GSM6160813 | Insulin | At-diagnosis | Female | 2 |
| ... | ... | ... | ... | ... |
| GSM6160833 | Placebo | Pre-clinical | Male | 2 |

s) ← ~ Treatment

# R: DE results

## Volcano Plot



- Each dot represents a gene;

- Most statistically significant genes will be towards the top;

- Most upregulated genes will be towards the right;

- Most downregulated genes are towards the left;

- Fold change equal to 0 means that the expression does not change.

# R: DESeq2

Let's get a count matrix from an RNA sequencing experiment from the National Center for Biotechnology Information repository!

Dataset: https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE232750

Home - GEO - NCBI (nih.gov)

GEO DataSets    [GEO DataSets ▾]  [_____]  [Search]
                        Advanced                                              Help

| DataSet | Title | Organism(s) | Platform | Series | ▶ Samples |
|---------|-------|-------------|----------|--------|-----------|
| GDS6063 | Influenza A effect on plasmacytoid dendritic cells | *Homo sapiens* | GPL10558 | GSE68849 | 10 |
| GDS6010 | Influenza virus H5N1 infection of U251 astrocyte cell line: time course | *Homo sapiens* | GPL6480 | GSE66597 | 18 |
| GDS5879 | Pulmonary CDC11c+ cells from young and middle-age animals | *Mus musculus* | GPL6885 | GSE71868 | 8 |
| GDS5826 | Multiple myeloma cell lines with acquired resistance to chemotherapeutic agent carfilzomib | *Homo sapiens* | GPL570 | GSE69078 | 12 |
| GDS5825 | Interleukin-1α deficiency effect on injured spinal cord | *Mus musculus* | GPL6246 | GSE70302 | 12 |
| GDS5881 | Nebulin deficiency effect on the soleus | *Mus musculus* | GPL6246 | GSE70213 | 12 |
| GDS5880 | Nebulin deficiency effect on the quadriceps | *Mus musculus* | GPL6246 | GSE70213 | 12 |
| GDS5913 | SRPIN803 small molecule inhibitor of SRPK1 effect on retinal pigment epithelial cell line | *Homo sapiens* | GPL570 | GSE62947 | 6 |
| GDS5665 | Pathogen-associated molecular-pattern curdlan effect on interleukin-2 deficient GM-CSF myeloid dendritic cells | *Mus musculus* | GPL6246 | GSE58120 | 12 |
| GDS5662 | Histone demethylase KDM3A-deficiency effect on estrogen-stimulated breast cancer cells in vitro | *Homo sapiens* | GPL10558 | GSE68918 | 11 |

About GEO DataSets          DataSet Browser                    SRA

Construct a Query           Programmatic Access

Download Options            GEO2R

# Time to do some coding!

## Let's go back to the R visualization script