# Introduction to basic R



Abu Bakar Siddique

Bioinformatician, SLUBI

# Contents

- About R

- R vs Rstudio

- Timeline

- Setting a project or working directory

- Interacting with R

- Packages

- Assign values or objects

- Data types

- Data structures

- How to import and export data or results

- Handle large data

- Housekeeping

- Outliers

# Basic about R

## R is ...

- a programming language

- a programming platform (= environment + interpreter)

- a software project driven by the core team and the community

- a very powerful tool for statistical computing

- a very powerful computational tool in general

## Yet ...

- it is very elegant

- it becomes more and more feature-rich

## R is not ...

- a tool to replace a statistician

- the very best programming language

- the most elegant programming solution

- the most efficient programming language

R: Engine | RStudio: Dashboard

## R

- Programming language

- For data analysis and graphics

- Refers to both the language and the software that interprets it's scripts

- Free and open source

## RStudio

- User interface for working with R

- Wrapper around the R language

- Extends what R can do and facilitates writing R code

- Free and open source

# Timeline

- ca. 1992 — conceived by [Robert Gentleman](#) and [Ross Ihaka](#) (R&R) at the University of Auckland, NZ as a tool for **teaching statistics**

- 1994 — initial version

- 2000 — stable version

- 2011 — RStudio, first release by J.J. Allaire

- ca. 2017 — Tidyverse by Hadley Wickham

# Setting up a project

# The working directory

- Where R will look for and save files

- Check working directory with the getwd() , setwd() functions

```
> getwd()
[1] "C:/Users/auue0001/OneDrive - Sveriges lantbruksuniversitet/Dokument"
```

```
> setwd("C:/Users/auue0001/OneDrive - Sveriges lantbruksuniversitet/Dokument")
```

# Suggested subdirectories



- data/             for raw data and intermediate datasets

- data_output/     modified versions of raw data

- documents/       outlines, drafts, other text

- fig_output/        graphics generated by scripts

- scripts/            R scripts for different analyses or plotting

```
# create subdirectory
dir.create(path = "path_to_working_directory/data_output")
```

# Interacting with R

# Packages

- developed by the community

- cover several very diverse areas of science/life

- uniformly structured and documented

- organized in repositories:
  - CRAN

# Packages



- developed by the community

- cover several very diverse areas

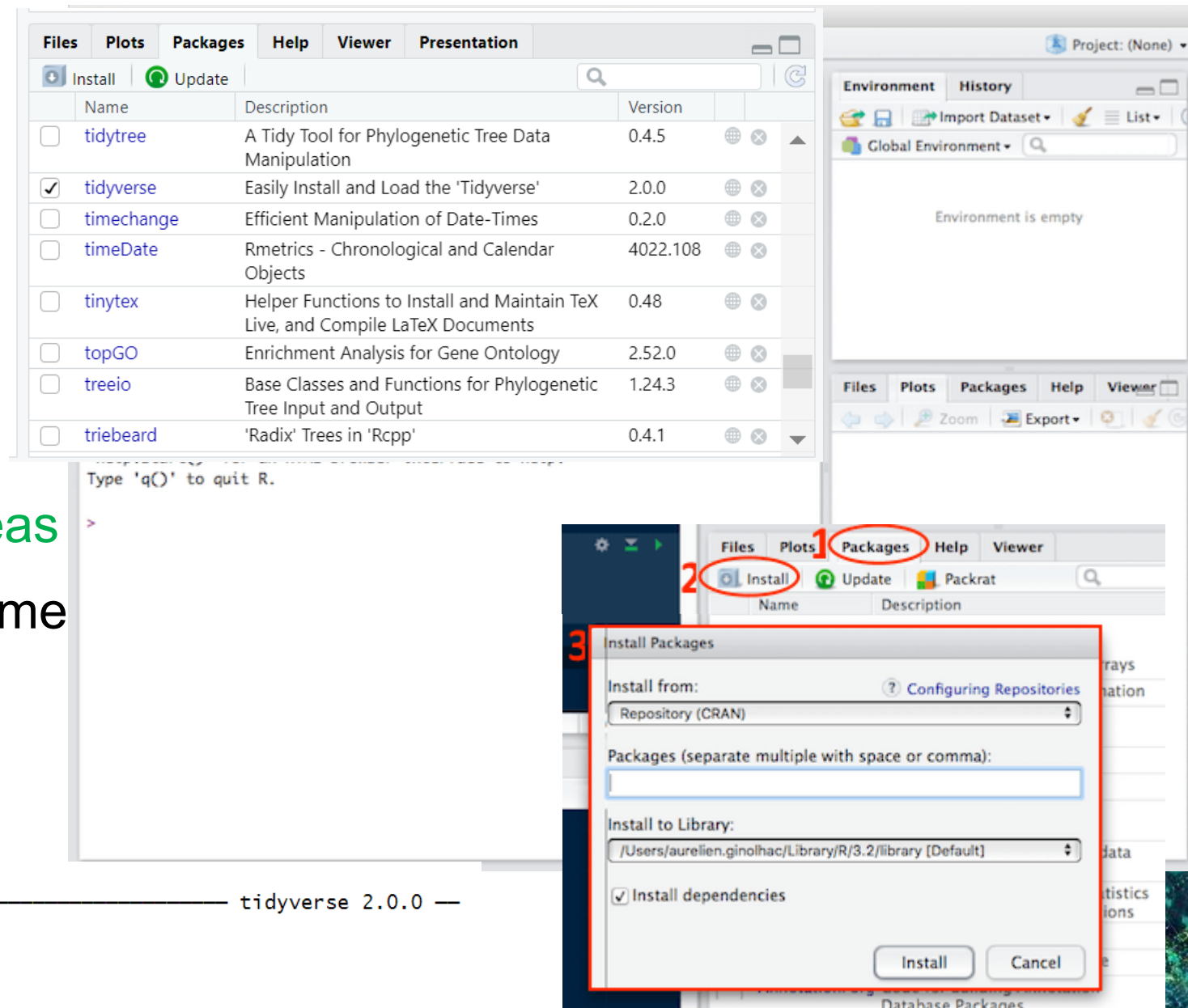- uniformly structured and docume

- organized in repositories:
  - CRAN

# Packages

- developed by the community

- cover several very diverse areas of science/life

- uniformly structured and documented

- organized in repositories:
  - CRAN
  - Bioconductor
  - R-Forge
  - GitHub

```
if (!require("BiocManager", quietly = TRUE))
    install.packages("BiocManager")
BiocManager::install(version = "3.18")
```

# Assign values to objects

# Assign values to objects

p <- 3


x <- c(1, 6, 8)


y <- c("car", "truck")

# Scalars, vectors and matrices

# Types of vectors and objects

- **Scalar** (0 dimention):

  p <- 3

- **Vector** (arrays -1 dimention) :

  x <- c(1, 6, 8)

  y <- ("car", "truck")

```
1  > vec1 = c(1,4,6,8,10)
2  > vec1
3  [1]  1  4  6  8 10
4  > vec1[5]
5  [1] 10
```

- **Matrix** (2 dimension):

```
1  mat=matrix(data=c(9,2,3,4,5,6),ncol=3)
2  > mat
3       [,1] [,2] [,3]
4  [1,]    9    3    5
5  [2,]    2    4    6
```

# Types of vectors and objects

- **character**: y <- ("car", "truck")

- **numeric**:  x <- c(1.1, 6.2, 8.4)

- **integer**: z <- 2

- **logical**: TRUE, FALSE

- **complex**: 1+4i (complex numbers with real and imaginary parts)

R has many functions to examine features of vectors and other objects:,

•class() - what kind of object is it (high-level)?

•typeof() - what is the object's data type (low-level)?

•length() - how long is it? What about two dimensional objects?

•attributes() - does it have any metadata?

# Structures

# Data structures

- Matrix

- Data frames

- Factors (r assign a level for each

  values)

- Arrays

- Lists

```
1  mat=matrix(data=c(9,2,3,4,5,6),ncol=3)
2  > mat
3        [,1] [,2] [,3]
4  [1,]    9    3    5
5  [2,]    2    4    6
```

```
1  > t = data.frame(x = c(11,12,14),
2   y = c(19,20,21), z = c(10,9,7))
3  > t
4      x  y  z
5  1 11 19 10
6  2 12 20 9
7  3 14 21 7
```

```
> x <- 1:12

> # Create a 3 x 4 array from the vector
> my_array <- array(values, dim = c(3, 4))

> # Print the array
> print(my_array)
     [,1] [,2] [,3] [,4]
[1,]   1    4    7   10
[2,]   2    5    8   11
[3,]   3    6    9   12
> values
 [1]  1  2  3  4  5  6  7  8  9 10 11 12
```

```
1  > L = list(one=1, two=c(1,2),
2   five=seq(0, 1, length=5))
3  > L
4  $one
5  [1] 1
6  $two
7  [1] 1 2
```

# Getting a dataset in R

# Ways to get data in R or in RStudio

- **Manually:**
  - data.frame() function in Base R,
  - or the tibble() function in the tidyverse.

- **Import it from a file:**
  - Fasta: VCF (vcfR package)
  - Text: TXT (readLines() function)
  - Tabular data: CSV, TSV (read.table() or read_delim() fu... readr package which contains read_csv())
  - Excel: XLSX (xlsx package)
  - Google sheets: (googlesheets package)
  - Statistics program: SPSS, SAS (haven package)
  - Databases: MySQL (RMySQL package)



Data import with the tidyverse :: CHEAT SHEET

```
Student ID,Full Name,favourite.food,mealPlan,AGE
1,Sunil Huffmann,Strawberry yoghurt,Lunch only,4
2,Barclay Lynn,French fries,Lunch only,5
3,Jayendra Lyne,N/A,Breakfast and lunch,7
4,Leon Rossini,Anchovies,Lunch only,
5,Chidiegwu Dunkel,Pizza,Breakfast and lunch,five
6,Güvenç Attila,Ice cream,Lunch only,6
```

Table 7.1 shows a representation of the same data as a table.

Table 7.1: Data from the students.csv file as a table.

| Student ID | Full Name | favourite.food | mealPlan | AGE |
|---|---|---|---|---|
| 1 | Sunil Huffmann | Strawberry yoghurt | Lunch only | 4 |
| 2 | Barclay Lynn | French fries | Lunch only | 5 |
| 3 | Jayendra Lyne | N/A | Breakfast and lunch | 7 |
| 4 | Leon Rossini | Anchovies | Lunch only | NA |
| 5 | Chidiegwu Dunkel | Pizza | Breakfast and lunch | five |
| 6 | Güvenç Attila | Ice cream | Lunch only | 6 |

variables

observations

values

# Processes



Import → Tidy → Transform

Visualize

Model

Understand

Program

# How to **export** or save as **results or data**

write.csv(df , file = "path/to/your/saving/folder/df.csv")


writexl::write_xlsx (test_df, path = "C:/Users/your_username/test_df.xlsx", col_names = TRUE, format_headers = TRUE)

# How to **export** or save as **plot**

```
dev.print (device=jpeg, file="path/to/your/saving/plots/figure_1a.jpg",
width=par("din")[1]*300, res=300, quality=100)


ggsave("path/to/your/saving/plots/figure_1a.png", plot = p1, bg ="white")
```

**!! Practical session (afternoon session): read a test csv file into Rstudio!!**

# Seeking help

- RStudio help interface
  - ?sum
  - help.search(), with term in "" inside parentheses

- Google or chatgpt "R <task>"

- When asking others
  - Use correct words
  - Reduce to reproducible example
  - Always include output of sessionInfo() function

Handle large datasets

Start by looking at the file names and sizes:

```r
fhvhv_csv_files <- list.files("original_csv", recursive=TRUE, full.names = TRUE)
data.frame(file = fhvhv_csv_files, size_Mb = file.size(fhvhv_csv_files) / 1024^2)
```

```
##                                                          file   size_Mb
## 1  original_csv/2020/01/fhvhv_tripdata_2020-01.csv 1243.4975
## 2  original_csv/2020/02/fhvhv_tripdata_2020-02.csv 1313.2442
## 3  original_csv/2020/03/fhvhv_tripdata_2020-03.csv  808.5597
## 4  original_csv/2020/04/fhvhv_tripdata_2020-04.csv  259.5806
## 5  original_csv/2020/05/fhvhv_tripdata_2020-05.csv  366.5430
## 6  original_csv/2020/06/fhvhv_tripdata_2020-06.csv  454.5977
## 7  original_csv/2020/07/fhvhv_tripdata_2020-07.csv  599.2560
## 8  original_csv/2020/08/fhvhv_tripdata_2020-08.csv  667.6880
## 9  original_csv/2020/09/fhvhv_tripdata_2020-09.csv  728.5463
## 10 original_csv/2020/10/fhvhv_tripdata_2020-10.csv  798.4743
## 11 original_csv/2020/11/fhvhv_tripdata_2020-11.csv  698.0638
## 12 original_csv/2020/12/fhvhv_tripdata_2020-12.csv  700.6804
```

We can already guess based on these file sizes that with only 4 Gb of RAM available we're going to have a problem.

# Handle large datasets

- Good management strategies for large files
    - if you work with **10 to100 GB** regularly!?
    - *r-datatable.com*
    - library(data.table)
    - R script:
        *firstscript.R*

Downsample in database → Pull Sample to Dev Machine → Build Model

Split data into logical chunks → Pull chunks in individually → Build separate models

Do compression operation in database → Pull only compressed data into R → Analytical work

# Good housekeeping strategies for scripts

- ## Comments, Structure:
  - – Use #
  - – Outline

- ## Consistent Naming Conventions:
  - – Use consistent naming conventions for variables, functions, and objects
  - – Avoid duplicating code. create a function or use a loop

- ## Version Control:
  - – Git. Platforms like GitHub or GitLab repositories.

- ## Imports and Dependencies:
  - – List all package imports at the top of your script
  - – call rhistory: sessionInfo()

- ## File Organization:
  - – Separate your R scripts, data, documentation, and output files into logical folders.

- ## Reproducibility:
  - – Use Quarto, R Markdown Documents:
  - – Documentation Files:
  - – Create README files

# Filtering of genotype data



- Subset data, metadata
  - Based on "rows (observations)" , "variable" or "values"
  - Filter

- Missing data ("NA", stat, summary)



height_sub_12 <- subset(data, height<12)

data_filter_na <- data[!is.na(data$height),]

- `mutate()` adds new variables that are functions of existing variables
- `select()` picks variables based on their names.
- `filter()` picks cases based on their values.
- `summarise()` reduces multiple values down to a single summary.
- `arrange()` changes the ordering of the rows.

# Conversion of different genotype file formats

Common programs used to handle the file formats:

- Hapmap: library(plink), library(VariantAnnotation)

- Numeric: library(readxl, writexl), library(jsonlite)
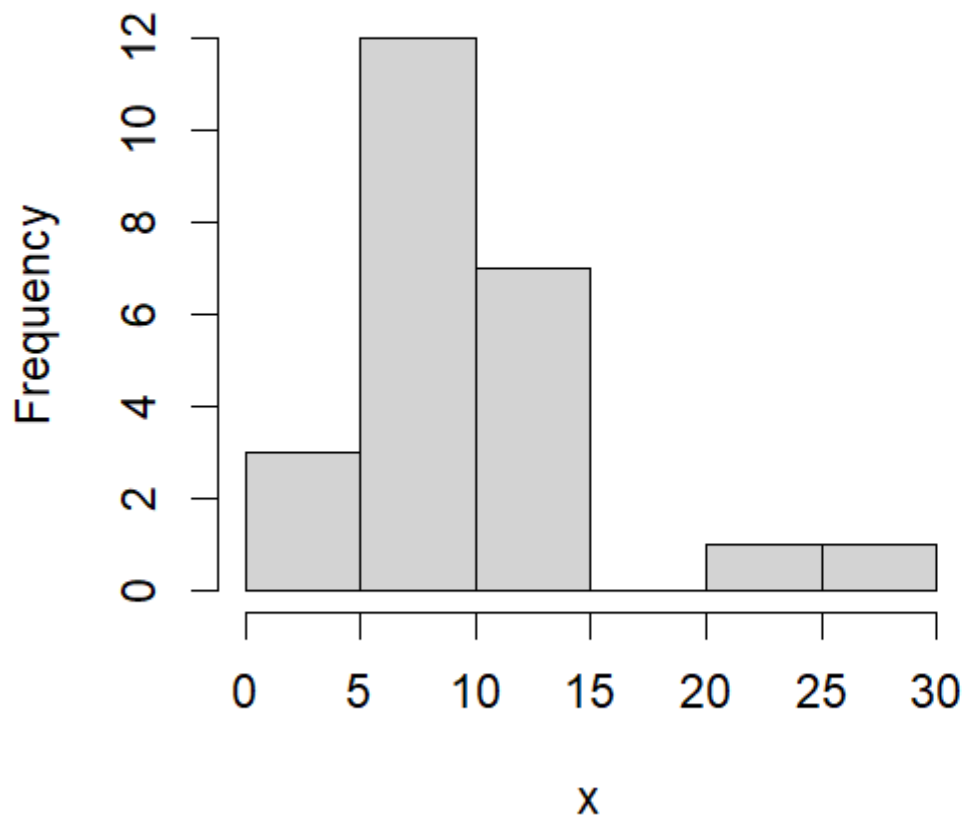
- Haploid format – one letter code (readLines function)

# Outliers in genotype and phenotype datasets

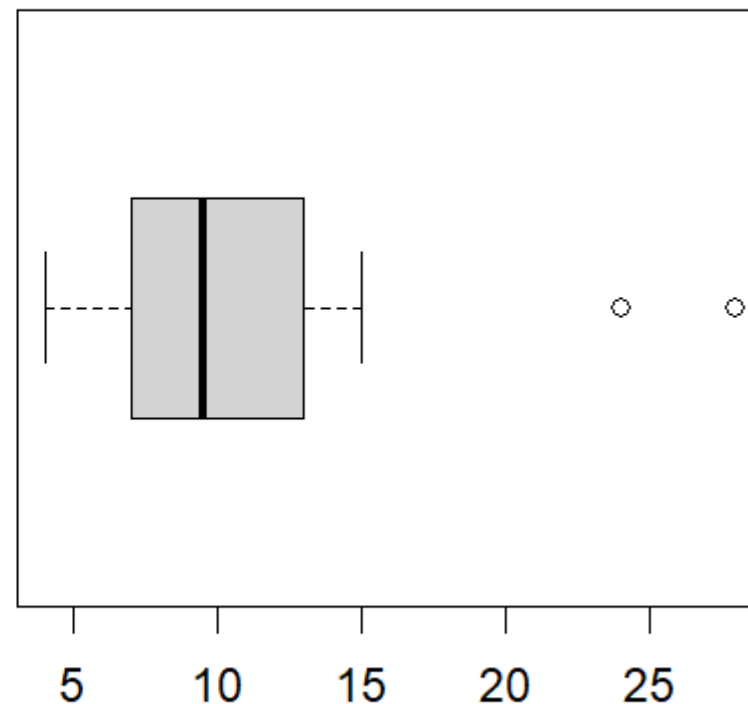- Histogram, scatterplot, and boxplot, Q-Q plot, chi square test

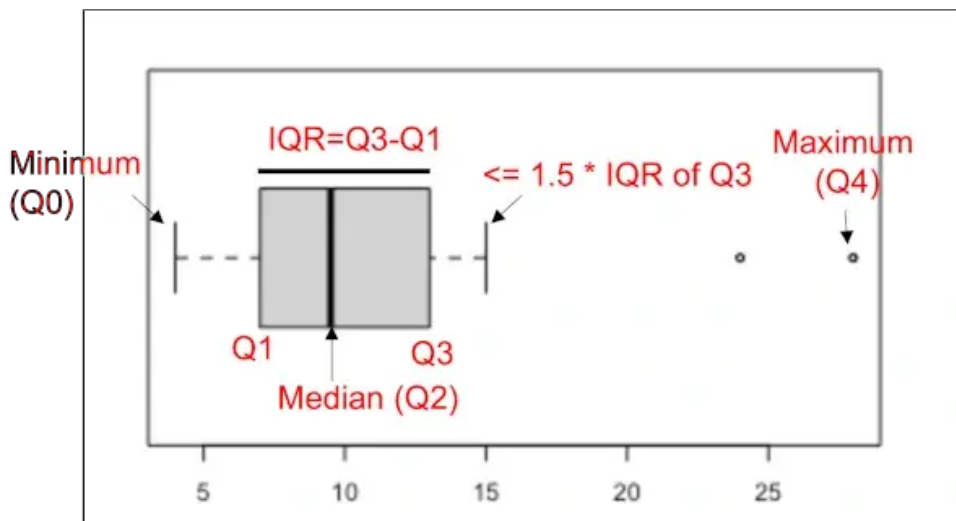x = c(10,4,6,8,9,8,7,6,12,14,11,9,8,4,5,10,14,12,15,7,10,14,24,28)

# Outliers in genotype and phenotype datasets

- Histogram, scatterplot, and boxplot, Q-Q plot, chi square test

- IQR



```
x = c(10,4,6,8,9,8,7,6,12,14,11,9,8,4,5,10,14,12,15,7,10,14,24,28)

# get values of Q1, Q3, and IQR
summary(x)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
  4.00    7.00    9.50   10.62   12.50   28.00

# get IQR
IQR(x)
[1] 5.5

# get threshold values for outliers
Tmin = 7-(1.5*5.5)
Tmax = 12.50+(1.5*5.5)

# find outlier
x[which(x < Tmin | x > Tmax)]
[1] 24 28

# remove outlier
x[which(x > Tmin & x < Tmax)]
[1] 10  4  6  8  9  8  7  6 12 14 11  9  8  4  5 10 14 12 15  7 10 14
```

# Questions?!

# Acknowledgements

- NBIS