

# Exercise day 1

## Introduction to R for Basic Statistics

Alessandra Meddis

### Exercise A:

For this exercise we will work with a subset of “follicle” data, collected from patients with cancer that had OTC (ovarian tissue cryopreservation). Follicles were cultured for 8 days and the diameter was collected every 2 days. The aim of the study was to compare the follicles growth among different treatment groups.

More information on the study and full data are available at <https://doi.org/10.1016/j.rbmo.2023.06.011>.

Information are collected in two data sets.

Data set *follicle* with:

- Patient: patient ID
- Number: follicle ID
- Day\_x: follicle diameter at Day x

Data set *patient* with:

- Patient ID
- Treatment group
- Type of Disease
- Age at Day 0

**Question 1** Load the data sets into R from the course material webpage by running:

```
dbf<-read.csv("https://raw.githubusercontent.com/AMeddis/IntrotoR-for-Basic-Statistics/main/data_exercise/follicle.csv")

dbp<-read.csv("https://raw.githubusercontent.com/AMeddis/IntrotoR-for-Basic-Statistics/main/data_exercise/patient.csv")
```

*Note: you can assign different names to the data by changing the left part of <-.*

Focus on the data set *patient*:

1. Check the dimension of the data.frame. How many patients were included in the study?
2. Visualize the first lines of the data using the function *head()*.
3. Print a summary of the data. What is the data type of each variable?

**Question 2** Treatment and Disease are *characters*, can we understand from the summary how many different diseases are in the data?

1. Would it be better if Treatment and Disease were of a different type? If yes, which one?
2. Transform them into factor. You can use the function *factor()* in R.
3. Print the summary of the data, can you see any difference?

**Question 3** Age is a continuous covariate. We would like to have an idea of the age distribution of patients included in the study:

1. Show min,max and mean for age
2. Calculate the mean of age for each disease group (you can use *aggregate* or *tapply*)
3. Create a new variable in the data.frame for a categorical variable for age with the median as cut-off for the two categories. We can code it using two different functions in **R**:
  - 3a. Use the function *cut* in **R**. Run the command *str(NameofDatabase)*, of which type is the new categorical variable?
  - 3b. Use the function *ifelse* in **R**. Run the command *str(NameofDatabase)*, of which type is the new categorical variable?

**Question 4** We now focus on disease and treatment groups:

1. Show the proportion of patients in each treatment group (use *prop.table*)
2. Show the number of patients in each disease group
3. Show the number of patients by treatment group and disease (two-ways table)
4. Create a new variable *Cancer* that groups Disease into *Breast cancer* and *Others* (you can use *ifelse()*)
5. Calculate the average age for each group respect to *Cancer*

**Question 5** Consider only patients with Breast cancer:

1. Subset data for Breast cancer patients (use the function *subset*)
2. Show the number of patients for each age group
3. Calculate mean and standard deviation for age