

# Introduction to R

## Session 03: Basics of programming

Álvaro Pérez<sup>1</sup>

Instituto Tecnológico Autónomo de México

Fall 2024

---

<sup>1</sup>Based on PhD Romero Londoño's notes.

- ▶ Basic idea: you have a vector of values, and you have an “iterator” variable
- ▶ You go through the vector, setting the iterator variable to each value one at a time
- ▶ Then you run a chunk of that code with the iterator variable set

```
1  for (iteratorvariable in vector) {  
2      code chunk  
3  }
```

```
1 data(mtcars)
2 abovemed <- mtcars %>% filter(cyl >= median(cyl))
3 belowmed <- mtcars %>% filter(cyl < median(cyl))
4 for (i in c('mpg', 'disp', 'hp', 'wt')) {
5   print(mean(abovemed[[i]]) - mean(belowmed[[i]]))
6 }
```

- ▶ Sometimes, you want R to do something only in some cases
- ▶ The basic syntax to do this is an “if” statement

```
1 if (test_expression) {  
2   statement  
3 }
```

- ▶ The “test\_expression” usually depends on logical operators
  - ▶ ‘&’ is AND
  - ▶ ‘|’ is OR
  - ▶ To check equality use ‘==’, not ‘=’
  - ▶ ‘>=’ is greater than OR equal to, similarly for ‘<=’

# For loops

- ▶ Sometimes is useful to “break” (stop) a loop, or skip an iteration
- ▶ Do this with “break” and “next”
- ▶ Often “break” and “next” are nested within an “if” statement

```
1 #printing odd numbers
2 m=20
3 for (k in 1:m){
4     if (!k %% 2)
5         next
6     print(k)
7 }
```

## Example: loop combined with if and break

```
1 #finding the first number divisible by 13 greater than
   100
2 for (k in 100:1000000){
3     if (!k %% 13){ #checks if number is odd
4         break
5     }
6 }
7 print(k)
```

# Example: finding prime numbers

```
1 for(num in 1:100){
2   # Program to check if the input number is prime or not
3   flag = 0
4   # prime numbers are greater than 1
5   if(num > 1) {
6     # check for factors
7     flag = 1
8     for(i in 2:(num-1)) {
9       if ((num %% i) == 0) {
10        flag = 0 #if number is divisible , then not prime
11        break #and we can break the loop
12      }
13    }
14  }
15  if(num == 2)    flag = 1
16  if(flag == 1) {
17    print(paste(num, " is a prime number"))
18  } else {
19    print(paste(num, " is not a prime number"))
20  }
21 }
```

- ▶ Functions are used when you have to repeat the same operation multiple times
- ▶ Avoids coding mistakes (and if there is one you only need to fix it once)
- ▶ Let's you break code into simpler parts which become easy to maintain and understand
- ▶ It's pretty straightforward to create your own function in R programming.

```
1 func_name <- function (argument) {  
2   statement  
3   return(whatyouwant)  
4 }
```



# Example: finding prime numbers

```
1 # Identifying prime number
2 prime_num<-function(num){
3   # Program to check if the input number is prime or not
4   flag = 0
5   # prime numbers are greater than 1
6   if(num > 1) {
7     # check for factors
8     flag = 1
9     for(i in 2:(num-1)) {
10       if ((num %% i) == 0) {
11         flag = 0 #if number is divisible , then not prime
12         break #and we can break the loop
13       }
14     }
15   }
16   if(num == 2)      flag = 1
17   return(flag)
18 }
```

# Example: finding prime numbers

```
1 prime_num(2)
2 prime_num(3)
3 prime_num(4)
4 prime_num(6131)
5 prime_num(4684561123)
```

- ▶ Functions can have multiple arguments.
- ▶ Functions can have default values for arguments.
- ▶ Functions can output vectors, list, or any other object.

# Example: power function

```
1 pow <- function(x, y = 2) {  
2   result1 <- x^y  
3   return(result1)  
4 }  
5 pow(3)  
6 pow(3,3)
```

- ▶ Automate everything that can be automated.
- ▶ Design your workflow so that when you add new data or tweak something you don't need to manually recreate anything
- ▶ Avoid copying and pasting figures/tables from one software into another
  - ▶ e.g., impossible to automate: copying a graph from R into Word
- ▶ General workflow: raw data → clean data → analysis → (pretty) table or figure

- ▶ Use the R studio browser or excel
- ▶ Scroll through the variables and accompany yourself with what you've got visually

- ▶ Check the size of your dataset using `dim`
- ▶ Check the structure of your data using `str`
- ▶ Check for missing variables
  - ▶ Sometimes the data has -99 for missing
  - ▶ Others is truly missing (i.e., NA)
  - ▶ Others is an empty string

# Always check your merges (combining data sets)

- ▶ During a stage of arranging datasets, you will likely merge
- ▶ Make sure you count before and after you merge so you can figure out what went wrong, if anything
- ▶ Also make sure you understand what type of merge you are using/need
  - ▶ many to one
  - ▶ one to many
  - ▶ one to one
  - ▶ many to many (please no!!!)



- ▶ Beyond checking the data, here are a few things that will prevent errors
  - 1. Organized subdirectories
  - 2. Automation
  - 3. Naming conventions
  - 4. Version control

- ▶ Things to keep in mind when naming:
  1. variables,
  2. datasets
  3. scripts
  
- ▶ Avoid
  - ▶ meaningless words (e.g., lmb2)
  - ▶ dating (e.g., temp05012020), except to save old version in the archive
  - ▶ numbering (e.g., outcome25)
  - ▶ any of these will confuse your future self

- ▶ Variables should be readable to a stranger
  - ▶ Say that you want to create the product of two variables. Use an underscore and mash the two together
    - ▶ `price_mpg <- price * mpg`
- ▶ Name the variable exactly what it is when possible
  - ▶ `bmi <- weight / (height*height * 703)`

- ▶ Always start your programming scripts with a header
- ▶ The header is a big comment on the code and says:
  - ▶ The purpose of the script
  - ▶ The author(s)
  - ▶ The date the script was created
  - ▶ The last date in which it was updated/modified (and by whom)

# Automating Tables and Figures (more on visualization below)

- ▶ Goal is to make “beautiful tables/figures” that are never edited post-production
- ▶ Tables/figures should be readable on their own (with the table/figure notes)
- ▶ Large fixed costs learning commands like “stargazer” in R or “estout” in R

- ▶ Maximize the data-to-ink ratio by using as little ink as possible to show your data
- ▶ Remove non data ink, e.g., extra gridlines
- ▶ Remove redundant data
- ▶ Remove indicators you don't need
- ▶ Display only as many decimal places as are relevant

# Choose the type of visualization based on the information you want to convey

- ▶ Use a table only if you can't use a figure. Usually when
  - ▶ You need to show exact numerical values
  - ▶ You want to allow for multiple localized comparisons
  - ▶ You have relatively few numbers to show
- ▶ Avoid pie charts
  - ▶ It is difficult for the audience to visually distinguish the size of each wedge
  - ▶ A good alternative for showing shares of a whole is a stacked bar chart

- ▶ Make the font larger than the Stata/R preset
- ▶ Label your axis
- ▶ The Y-axis tickmarks ideally should be horizontal (R puts them sideways)
- ▶ Use legends if you have multiple lines
- ▶ Avoid having two-axis
- ▶ Avoid pie-charts



- ▶ Add x-axis tick marks
- ▶ We only want people aged 5-18
- ▶ Make the y-axis in percent terms
- ▶ Use a more appealing color scheme
- ▶ The y-axis is misleading, lets make it go to 100
- ▶ Add a legend to explain what the colors mean
- ▶ Make the numbers in the y-axis easier to read and remove legend box