# Introduction to R programming

Iván Andrés Trujilllo Abella

ai-page.readthedocs.ios

**ivantrujillo1229@gmail.com**

# History of R

R is a dialect (an implemntation) of S language.

> ...
>
> Created in 1991 in new zealand.

# fragile

```
1 # The print function
2 print('hello world')
```

Assigment operation,

```
1 price <- 30 # assigment operator
2 quantity <- 5
3 income <- price * quantity
```

### vectorized

R objects as used before are really vectors.

```
output_msg <- "profit"
```

**...**

The object is character vector in R, the first is profit

[1] 10 a vector whose first element is the number ten
We'll create a sequence with colon(:) operator.

```
1 ito <- 1:10
```

# R objects

atomic classes of objects:

- character
- integer
- numeric (float)
- complex
- logical ( booleans)

# keep information

**vector - list**

### vector

each one of its elements are of the same class

### list

can keep information of all kind

```
vector() # i don't have idea then;
help("vector")
```

# Numbers

in R are generally treated as numeric objects ( double precision real numbers).

```
1 123L # explitcit integer
```

### Special
- inf (infinity)
- NaN (Not A Number)

# Attributes

names, dimensions, class, lenght and another.

```
1 help("attributes")
2 attributes(df) # rows and columns
```

# Vector - list

the $c()$ (concatenate) function crate **vectors** of objects.

```r
df <- c(0.9, 0.87, 0.95) # numeric object
booleans <- c(TRUE, FALSE, T, F) # logical
```

### coercion

if you mix different kind of objects in a vector then occur coercion casting all the object to the same class, for instance c('1', 2) caste 2 to '2'.

# Type

In programming languages define the type of data is very important thus:

```
1 bool <- c(True, F)
2 class(bool)
3 # explicit conversion
4 as.numerical(bool)
5 as.character(bool)
6 binaries <- c(1,0,1)
7 as.logical(binaries)
8 vowels <- c('a','e','i')
9 as.logical(vowels)
```

```
1 data <- list('True', FALSE, 0 , 1.0)
```

Then the first element is a vector of characters, then second is a vector (logical) with the element FALSE, etc...

# Matrices

```
1 mat <- matrix(nrow=3, col= 2)
2 attributes(mat)
3 dim(mat)
```

```
1 matrix(1:4,nrow=2,ncol=2) # column-wise (left-upper
    )
2 data <- 1:4
3 dim(data) <- c(2,2)
4 print(data)
```

# cbind and rbind

```r
x <- 5:10
y <- 50:55
cbind(x,y)
rbind(x,y)
```

# Factors

Useful in categorical data, ordinal or cardinal, think in the problem itself
that carry out pass categorical to numbers low=1, medium =2, high=3.

```
1 var <- factor(c('high', 'medium','low', 'low'))
2 table(var)
3 unclass(var) # how the factors are represented by R
```

# Baseline factor

```
var <- factor(c('cs','math','cs' ,'economics',
'math','cs'),
leves=c('math', 'cs', 'economics'))
```

The first level is used as baseline level in some statistical procedures.

# Missing values
**NA and NaN**

- NA (Not Avaliable) is general, a dont have information about missingness is a class!.
- NaN (Not a Number) arise from undefined mathematical operations as $\frac{a}{0}$

```
is.nan()
is.na()
```

a NaN (related to numerical) is also a NA but not viceversa.

The following is very important handled missing values.

```
1 data <- c (1 ,2 ,4 , NA , NaN )
2 is . nan ( data )
3 is . na ( data
4 data <- c ( '1 ' , 2 , NaN , NA )
5 is . nan ( data )
6 is . na ( data )
```

# DataFrame

Tabular data, important points:

- unlike matrices can handled different types of data (like list)
- could be comverted to matrices as df.matrix()

Create a dataframe to test:

```
df <- data.frame(var1=c('red','red','blue', 'blue')
    ,
var2=c(T, T, F, F), var3=c(1,1,0,0))
# Describe the data
nrow(df)
ncol(df)
```

# Names

```r
booleans <- c(T,F,T,F)
#Give name to each element...
names(booleans) <- c("yes", 'not', 'si', 'no')
names(booleans)
```

# name and indexing

```
1 data <- list(age= 40, id = 0, ibm =23)
2 data$age
3
4 mat <- matrix(1:6, ncol =3, nrow=2)
5 dimnames(mat)<-list(c('A', 'B'),
6 c('col1', 'col2', 'col3'))
7 # try
8 rownames(mat) <- c("A", "B")
9 colnames(mat) <- c('col1', 'col2', 'col3')
10 # indexing
11 mat['A', 'col3']
```

# read and write common

- read.csv()
- read.table()
- wite.table()
- save()

Remeber that is very important get information in documentation
`help("read.table")` also `read.csv()` is identical to `read.table` but if
separator by deafult is the comma(,) character. important also remark
that R skip lines that begin with # character.

# subsetting

- [ hold the class
- [[ dont hold the class (think in lists)
- $ extract elemens (named) in list or dataframe

```
1 data <- c('Blue', 'Green', 'Red')
2 data[1] # take the first element (numerical
     index)
3 data[1:2]# from first to second
4 data[c(T,F,T)] #logical index
5 data[data=='Green'] # two equals signs (=).
6 data[data=='Green' | data=='Blue']
7 data <- c(5,10,11,2, 7)
8 data[data>6]
```

# Research

...

- Research about the function `paste`
- Research about the function `cat`

# subsetting

**list**

```
1 data <- list(var1= 5:10, var2 = 15:20)
2 data[1] + 1   # error is a list!
3 data[[1]] + 1   # numeric vector
4 class(data$var1)
5 class(data['var1'])
6 class(data[['var1']])
```

is important take care, given that could be produce errors!! Single Square Bracket always return an object of the same class of the original object. Double Square bracket with name is equal to dollar sign.

# extract multiple elements

Uses single bracket operator

```
data <- list(a=10:30)
data$b <- 50:80
data$c <- 100:130
data[c("a","b")] # try with numerical index
# dollar sign not allow computed index for instance
    ;
column <- 'b'
data$column # return NULL
data[[column]] #execute!
data$b # execute!
```

subsetting nested elements.

```
data <- list(col1=list(T,F,T),
col2=list('Red','Blue'))
data[[c(1,2)]]
data[[1]][[3]]
```

# Subsetting matrix

$M_{ij}$ = M[i,j] where $i$ is rows and $j$ columnns.

```r
mat <- matrix(1:10, 2,5)
mat[1,4]
mat[1,]
mat[,4]
class(mat[1,3])
class(mat[1,3, drop=F])
# A rule is broken [] only return the same class
class(mat[1,])
class(mat[1,,drop=F])
```

# Removing missing values

```
1 bools <- c(T,F)
2 !bools
3 data <- c(1, 0 , 0.5, NaN, NA)
4 data[!is.na(data)]
5 x <- c(NA, 3, NaN, 5, 10)
6 y <- c(10, 11, 4, Nan, 4)
7 complete.cases(x,y)
8 # assign to both arrays to extract the complete
    data.
```

# in DataFrame

```r
data <- data.frame(
age = c(32, NaN, 31, NaN),
score = c(NaN, 70, 81, 90))
index <- complete.cases(data)
data[index,]
```

# Vectorized operations

```
1 x <- 10:15
2 y <- 0:5
3 x < y
4 x + y
5 x * y
6 x / y
```

```r
mat <- matrix(1:9, 3,3)
mat2 <- matrix(5:13, 3, 3)
mat*mat2 # Element wise multplication
mat %*% mat2 #Matrix multiplication
```

# Control structures

```
1 if (condition)
2 {
3 statement 1
4 } else{
5 statement 2 (when cond is FALSE)
6 }
```

# Example

```
1 price <- 30
2 quantity <- 50
3 income <- price * quantity
4 if (income >10){
5   print("Excellent!")
6 } else{
7   print(':(')
8 }
```

# else if

```
1 if (condition){
2   Statement 1
3 } else if (condition) {
4 Statement alternative
5 } else{
6 else statement
7 }
```

**...**

```r
price <- 10
quantity <- 1
income <- price * quantity
if (income >10){
  print("Excellent!")
} else if (income == 10){
  print('equlibrium')
} else {
  print(':(')
}
```

# loop

**for**

```
for (num in 1:10){
  print(num)
}
```

# Examples

```
1  x <- c(10, 0, 5)
2  for (i in x){
3    print(i)
4  }
5
6  for (i in seq_along(x) ){
7    print(i)
8  }
9
10 for (i in 1:length(x)){
11   print(x[i] + 1)
12 }
```

# Nested loop

```r
for (i in 1:2) {
for (j in 3:5){
  print(c(i,j))
}}
```

As excercise nested a matrix! ncol and nrow.

# Example

...

See simulation mean

See simulation estimator counter

See simulation

# while loop

```r
acum <- 1
while (acum < 5) {
print(acum)
acum <- acum + 1
}
```

# OR and AND

- && logical AND operator
- & element wise logical operator used in vectors to test AND
- || logical OR operator
- | element wise logical operator used in vectors to test OR

```
1  T && F
2  T || F
3  x <- c(T, T, F, F)
4  y <- c(T, F, T, F)
5  x&y
6  x|y
```

```r
# Search a number in the vector!
array <- c(10 ,40 , 50 , 1)
flag = FALSE
to_search = 1
index <- 1
while (index <= length(array) && flag ==FALSE){
  print(array[[index]])
  if (array[[index]] == to_search){
     flag <- TRUE
  }
  index <- index + 1
}
flag
```

The conditions are evaluated from left to right.

# Repeat, Next, Break

This topics are to research...

- `repeat` is a infinite loop whose exit condition is `break` statement.
- `next` skip a loop iteration (make an example with NAN)

# Functions

```r
function_name <- function(arguments){
function statement
}
```

```r
future_val <- function(initial_val, r, n){
initial_val*(1 + r)^n #you could uses return
}
```

# Excercise

...

write a function that compute the number of periods ($n$)

# Excercise

write a function that compute the number of periods (*n*)

```
periods <- function(lambda, r){
return(log(lambda)/log(r +1 ))
}
```

## Pi simulation

(Click here)

# Practical

...
See lab

# notes

Default arguments works like in python the `return` statement nos is necessary R return the last statement in the function body. THe arguments work as in python by name or position.
always check the arguments

---

...

args(print)

Important **lazy evaluation** evaluate arguments if they are needed.

```
# lazy evaluation
power <- function(a,b,c){
  a^(b)
}
power(10,2)
#   error   not is arised
```

# Important things in functions

# ... argument (ellipsis)

is important be careful using ellipsis (a undefined number of parameters).

```
suma <- function (...) {
args <- list (...)
acum <- 0
for (arg in args){
acum <- acum + arg}
acum
}
suma (1,2,3,4,5,6,7,8,9,10)
print(10*11/2)
```

ellipsis working as wrapped!

```
1  mean <- function (... , message ='the mean is '){
2    print ( message )
3    suma (...) / length ( list (...) )
4  }
5  mean (1 ,3 ,4) # arguments appear after of ellipses
      then they must be named explicity .
```
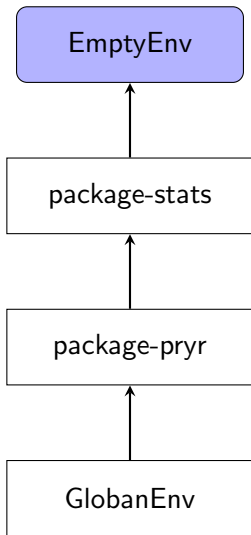
# Scoping rules

The scope is the place where a variables is defined and its visibility.

- Lexical scoping
- Dynamic scoping

in R are exist *free variables*.

Figure: Enviroment tree

First the symbol: value is saerched in the global enviroment, and after in each on of the packages listed in the following command.

```
search()
```

This load all packages that are loaded currently in R. Global enviroment is always de fisrt and base the last, when you load a package the namespace occupy the second place and all others shifdown a position.

- R functions are stored and searched in the same way of other object, therefore are searched in the environment tree.

# Annonymus functions

```
1  sapply ( c (1 ,2 ,3 ,4) , function ( x )  x ^2)
```

# $<< -$ **operator**

assign the variable to the global environment (not local where was defined).

```
1  local_var <- function(a){
2      result <- a^2
3      result
4  }
5  local_var(10)
6  result # object 'result' not found!
7
8  global_var <- function(a){
9  result <<- a^2
10  result
11  }
12  global_var(10)
13  result
```

# Closures

Analyze

```
1  counter <- function(){
2  counter <- 0
3  function(){
4  counter <<- counter + 1
5  counter
6  }
7  }
8  acum <- counter()
9  acum()
10 acum()
11 acum()
```

# Closures

- A function defined inside of another capture information from the outer function (its environment)
- The inner function can access and modify the variable count even though the function has already finished execution.

# closures

```
1  experiment <- function(sample_size){
2  mean <- 0
3  function(){
4  mean <<- mean(rnorm(sample_size))
5  mean
6  }
7  }
8  low_size <- experiment(10)
9  medium_size <- experiment(35)
10 for (i in 1:5){
11 print(c(low_size(), medium_size()))
12 }
```

Extend the function! to confidence intervals.

# Closure

```r
pi_estimator <- function(trials){
result <- function(){
acum <- 0
for (i in 1:trials){
if (sum(runif(2,-1,1)**2) <= 1){
acum <- acum + 1
}}
acum/trials*4
}}
small_trials <- pi_estimator(1000)
big_trials <- pi_estimator(1000000)
print(c(small_trials(), big_trials()))
```

# Closures

Given lexical scoping low_size and medium_size can reference its
environment definition!! and therefore the sample size defined 10 and 35
respectively.

```
ls(enviroment(low_size))
get('trials', environment(small_trials)) # knowing
    its enviroment!
```

# Exercise

```
1 var <- 1
2 call <- function() var
3 change <- function() {
4    var <- var +  100
5    #var <<- var +  100
6   call()
7 }
8 change() # the value is?
```

# Important point!

define and after call a function in global environment resemble to dynamic scoping given that calling and defining environment are the same (add an example).
Analyze

```
1 y <- 3
2 f <- function(x){
3 y <- 2
4 y^2 + eval(expression(g(x)),envir = parent.frame())
      # comment this line
5 #y^2 + g(x)  # uncomment this line
6 }
7 g <- function(x){
8 x^y
9 }
10 f(2)
```

# Dates and time

- Dates are stores as the number of days since 01-01-1970
- Times as the number of seconds from 01-01-1970

```
date_ <- as.Date('1971-01-01')
unclass(date_)
```

# Times

represented by POSIXct or POSIXlt

```
1 now <- as.POSIXlt(Sys.Date())
2 names(unclass(now))
3 now$mday
4 Sys.Date() # is in POSIXct format
5 # The number of seconds since 01-01-1970.
```

# strptime function

Data is written in a different format

```
date_str <- c("December 29, 2024 9:10")
date_ <- strptime(date_str, "%B %d, %Y %H:%M")
date_
start <- as.Date("2021-01-03")
finish <- as.Date("2022-02-01")
finish - start # take care with the compatibility
```

# Research time

- apply
- sapply
- lapply
- tapply
- mapply

# Excersice Write standard deviation function

# Excersice Write standard deviation function

```
 1 std_ <- function (..., ddof =0){
 2 args <- as.numeric(list(...))
 3 acum <- 0
 4 mean_ <- mean(args)
 5 for (arg in args){
 6   acum <- acum + (arg - mean_)**2
 7 }
 8 (acum / (length(args) - ddof))^0.5
 9 }
10 print(std_(10,3,1,4, ddof=1))
11 print(sd(c(10,3,1,4)))
```

# lapply

As an argument, receive a list (or cast it to one) and a function that receive as input each element of the list.

```
1 lapply (1:4 , function (x){x^2})
```

```
1 lapply # ellipses represent arguments in the
    function .
2 data <- list(list (10 ,3 ,1) ,
3 list (10 ,3 ,1 ,4) ,
4 list (10 ,3))
5 lapply (data , function (x){sd(unlist (x))})
6 lapply (data , function (x){do.call (std_ , x)}) #unpack
7 lapply (data , function (x) { do.call (std_ ,
8 c(x, ddof = 1)) }
```

List in R are "similar" to dictionary

```r
lapply(list(col1=1:5, col2=10:15),
mean, trim=0, na.rm=T) # note that are passed  mean
    function arguments!
```

# sapply

simplify the result of `lapply`

```
lapply(1:4, function(x) x**2)
sapply(1:4, function(x) x**2) # a vector is
    returned!
```

# apply

apply usually used to apply a function over a row or column of a matrix. The first argument is data, the second is the axis(1 or 2) and the third the function.

```
mat <- matrix(c(0, 1, 10,
                1, 0, 3,
                5, 0, 0),
              nrow = 3, byrow = TRUE)
apply(mat, 1, sum) # sum all values in each row
apply(mat, 2, sum) # sum all values in each column
```

# Research

- rowSums
- colMeans
- ...

how works

- quantile

- mapply
- tapply

# split

Split objects into groups determined by factors.

```r
data <- c(rnorm(5, 1), rnorm(5, 45), rnorm(5, 1000)
    )
f <- gl(3, 5) # three groups with five observations
     each one
split(data, f) # three groups
lapply(split(data,f), mean)
```

# apply - split

```
df <- data.frame(income=c(10, 14, 15, 23.2, 21.3),
exp = c(6,8,3, 2, 4),
group=c('math','math','math','cs','cs'))
split(df, df$group)
lapply(split(df, df$group), function(x) colMeans(x
    [, c('income', 'exp')]))
```

# more than one factor

```r
df <- data.frame(income=c(10, 14, 15, 23.2, 21.3,
    21.4),
exp = c(6,8,3, 2, 4,5),
group=c('math','math','math','cs','cs', 'cs'),
sex = c('female', 'male', 'male', 'female', 'female
    ', 'male'))

lapply(split(df, interaction(df$group, df$sex)),
    function(x) colMeans(x[, c('income', 'exp')]))
lapply(split(df, list(df$group, df$sex)), function(
    x) colMeans(x[, c('income', 'exp')])) #split
    allow us uses list
```

# Debug

research about

- `traceback`
- `debug`
- `browser`
- `trace`
- `recover`

# str

str() help us to visualize the data!

```
data <- list(1,2,3,4,5,6,7,8)
str(data)
str(df)
```

# Simulation

```r
set.seed(1) # uses seeds!
par(mar = c(4, 4, 1, 1))
domain <- seq(-3, 3 , length.out = 100)
range <- dnorm(domain)
plot(domain, range, type='l')
cumulative_range <- pnorm(domain)
plot(domain, cumulative_range, type='l')
```

# Research statistica distributions

**...**

```
help("distribution")
```

- rpois
- dt
- dnbinom
- ...

**prefixes**

- d(ensity)
- r(andom number generator)
- p(robability - cumulative distribution)
- q(uantile function)

# linear reg

```
1  linear <- function(b0, b1, x, u){
2  b0 + b1*x + u
3  }
4  sample_size <- 500
5  b0 = 1
6  b1 = 3.4
7  x <- rnorm(sample_size, 10, 3)
8  u <- rnorm(sample_size, 0, 1)
9  y <- linear(b0, b1, x, u)
10 lm(y ~ x )
```

# Sample

```
1 data <- c(100, 20, 5, 71, 89, 32, 121)
2 sample(data,2)
3 sample(data, 2, replace=T)
```