# Introduction to R*

## Part 2: Atomic Data Types - Homogeneous vectors

### Wim R.M. Cardoen

### Last updated: 07/25/2023 @ 13:07:10

## Contents

---

R can be summarized in **three** principles (John M. Chambers, 2016)

- Everything that exists in R is an **object**.
- Everything that happens in R is a **function** call.
- **Interfaces** to other languages are a part of R.

# 1 R Objects

- An object in R is (internally) represented as a pair: (**symbol**, **value**).
- A **symbol** is assigned a **value** by the use of an arrow pointing to the left (**<-**).
- There are **less favored** ways:
    - A simple equality sign (**=**).
    - Using the **assign()** function.

## 1.1 Examples

- Clean up the global environment i.e. remove all objects from the current R environment. Recommended!

```
rm(list=ls())
ls()
```

```
character(0)
```

- **preferred** way to assign variables

```
x <- 5.0
x
```

```
[1] 5
```

- alternative 1: mainly used to assign default function arguments

```
y = 5.0
y
```

```
[1] 5
```

```
mysamplevariance <- function(x, av=0){

  n <- length(x)
  if(n>1){
      return(1.0/(n-1)*sum((x-av)^2))
  }
  else{
      stop("ERROR:: Dividing by zero (n==1) || (n==0) ")
  }
}

x <- rnorm(10)
mysamplevariance(x)
```

```
[1] 0.3696939
```

```
mysamplevariance(x,mean(x))
```

```
[1] 0.3185038
```

```
var(x)
```

```
[1] 0.3185038
```

- alternative 2: even less used

```
assign("z", 5.0)
z
```

```
[1] 5
```

- functions are objects

```
f <- mean
f
```

```
function (x, ...)
UseMethod("mean")
<bytecode: 0x55a25d7957b0>
<environment: namespace:base>
```

```
val <- f(1:10)
val
```

```
[1] 5.5
```

"Nothing exists except atoms and empty space; everything else is opinion".  (Democritos)

# 2   Atomic Data Types

## 2.1   The core/atomic data types

- R has the following 6 **atomic** data types:
    - logical (i.e. boolean)
    - integer
    - double
    - character (i.e. string)
    - complex
    - raw (i.e. byte)

The latter 2 types (i.e. complex and especially raw) are less common.

The **typeof()** function determines the **INTERNAL** storage/type of an R object.

### 2.1.1   Examples

- boolean/logical values: either **TRUE** or **FALSE**

```
x1 <- TRUE
x1
```

```
[1] TRUE
```

```
typeof(x1)
```

```
[1] "logical"
```

- integer values ($\in \mathbb{Z}$):

```
x2 <- 3L
x2
```

```
[1] 3
```

```
typeof(x2)
```

```
[1] "integer"
```

- double (precision) values:

```
x3 <- 3.14
x3
```

```
[1] 3.14
```

```
typeof(x3)
```

```
[1] "double"
```

- character values/strings

```r
x4 <- "Hello world"
x4
```

```
[1] "Hello world"
```

```r
typeof(x4)
```

```
[1] "character"
```

- complex values ($\in \mathbb{C}$):

```r
x5 <- 2.0 + 3i
x5
```

```
[1] 2+3i
```

```r
typeof(x5)
```

```
[1] "complex"
```

## 2.2 Operations on atomic data types

- **logical** operators: **==**, **!=**, **&&**, **||**, **!**
- **numerical** operators: **+**, **-**, **\***, **/**, **^**, **\*\*** (same as the caret), but also:
  - integer division: **%/%**
  - modulo operation: **%%**
  - **Note**: matrix multiplication will be performed using **%\*%**
- **character/string** manipulation:
  - **nchar()**:
  - **paste()**:
  - **cat()**:
  - **sprintf()**:
  - **substr()**:
  - **strsplit()**:
  - **Note**: Specialized R libraries were developed to manipulate strings e.g. *stringr*
- explicit **cast**/conversion: https://data-flair.training/blogs/r-string-manipulation/
  - **as.{logical, integer, double, complex, character}()**
- explicit **test** of the type of a variable:
  - **is.{logical, integer, double, complex, character}()**

### 2.2.1 Examples

- Logical operators:

```
x <-3
y <-7
(x<=3) &&(y==7)
```

```
[1] TRUE
```

```
!(y<7)
```

```
[1] TRUE
```

- Mathematical operations

```
2**4
```

```
[1] 16
```

```
7%%4
```

```
[1] 3
```

```
7/4
```

```
[1] 1.75
```

```
7%/%4
```

```
[1] 1
```

- String operations

```r
s <- "Hello"
nchar(s)
```

```
[1] 5
```

```r
news <- paste(s,"World")
news
```

```
[1] "Hello World"
```

```r
sprintf("My new string:%20s\n", news)
```

```
[1] "My new string:         Hello World\n"
```

```r
city <- "Witwatersrand"
substr(city,4,8)
```

```
[1] "water"
```

- Conversion and testing of types

```r
s <- "Hello World"
is.character(s)
```

```
[1] TRUE
```

```r
s1 <- "-500"
is.character(s1)
```

```
[1] TRUE
```

```r
s2 <- as.double(s1)
is.character(s2)
```

```
[1] FALSE
```

```r
is.double(s2)
```

```
[1] TRUE
```

```r
s3 <- as.complex(s2)
s3
```

```
[1] -500+0i
```

```r
sqrt(s3)
```

```
[1] 0+22.36068i
```

## 2.3 Exercises

- 
  - Calculate $\log_2(10)$ using R's **log()** function
  - Perform the inverse operation and check that you get 10 back

- Let $z = 3 + 4i$

  - Use R's **Re()**, **Im()** functions to extract the real and imaginary parts of z.
  - Calculate the modulus of $z$ using R's **Mod()** function and check whether you the same answer using $\sqrt{\Re(z)^2 + \Im(z)^2}$.
  - Calculate the argument of $z$ using R's **Arg()** function and check whether you have the same answer using $\arctan\left(\frac{\Im(z)}{\Re(z)}\right)$.

# 3  Atomic vectors

- An **atomic** vector is a data structure containing elements of **only one atomic** data type. Therefore, an atomic vector is **homogeneous**.
- Atomic vectors are stored in a **linear** fashion.
- R does **NOT** have scalars:
  - An atomic vector of **length 1** plays the role of a scalar.
  - Vectors of **length 0** also exist (and they have some use!).
- A **list** is a vector not necessarily of the atomic type.
  A list is also known as a **recursive/generic** vector (*vide infra*).

## 3.1  Creation of atomic vectors

Atomic vectors can be created in a multiple ways:

- Use of the **vector()** function.
- Use of the **c()** function (**c** stands for concatenate).
- Use of the column operator **:**
- Use of the **seq()** and **rep()** functions.

The length of a vector can be retrieved using the **length()** function.

### 3.1.1  Examples

- use of the **vector()** function:

```r
x <- vector()  # Empty vector (Default:'logical')
x
```

```
logical(0)
```

```r
length(x)
```

```
[1] 0
```

```r
typeof(x)
```

```
[1] "logical"
```

```r
x <- vector(mode="complex", length=4)
x
```

```
[1] 0+0i 0+0i 0+0i 0+0i
```

```r
length(x)
```

```
[1] 4
```

```r
x
```

```
[1] 0+0i 0+0i 0+0i 0+0i
```

```r
x[1] <- 4
x
```

```
[1] 4+0i 0+0i 0+0i 0+0i
```

- use of the **c()** function:

```
x1 <- c(3, 2, 5.2, 7)
x1
```

```
[1] 3.0 2.0 5.2 7.0
```

```
x2 <- c(8, 12, 13)
x2
```

```
[1]  8 12 13
```

```
x3 <- c(x2, x1)
x3
```

```
[1]  8.0 12.0 13.0  3.0  2.0  5.2  7.0
```

```
x4 <- c(FALSE,TRUE,FALSE)
x4
```

```
[1] FALSE  TRUE FALSE
```

```
x5 <- c("Hello", "Salt", "Lake", "City")
x5
```

```
[1] "Hello" "Salt"  "Lake"  "City"
```

- use of the column operator:

```
y1 <- 1:10
y1
```

```
 [1]  1  2  3  4  5  6  7  8  9 10
```

```
y2 <- 5:-5
y2
```

```
 [1]  5  4  3  2  1  0 -1 -2 -3 -4 -5
```

```
y3 <- 2.3:10
y3
```

```
[1] 2.3 3.3 4.3 5.3 6.3 7.3 8.3 9.3
```

```
y4 <- 2.0*(7:1)
y4
```

```
[1] 14 12 10  8  6  4  2
```

```
y5 <- (1:7) - 1
y5
```

```
[1] 0 1 2 3 4 5 6
```

- **seq()** and **rep()** functions

```
z1 <- seq(from=1, to=15, by=3)
z1
```

```
[1]  1  4  7 10 13
```

```r
z2 <- seq(from=-2,to=5,length=4)
z2
```

```
[1] -2.0000000  0.3333333  2.6666667  5.0000000
```

```r
z3 <- rep(c(3,2,4), time=2)
z3
```

```
[1] 3 2 4 3 2 4
```

```r
z4 <- rep(c(3,2,4), each=3)
z4
```

```
[1] 3 3 3 2 2 2 4 4 4
```

```r
z5 <- rep(c(1,7), each=2, time=3)
z5
```

```
[1] 1 1 7 7 1 1 7 7 1 1 7 7
```

```r
length(z5)
```

```
[1] 12
```

### 3.1.2 Exercises

- Use the **seq()** function to generate the following sequence:
  6 13 20 27 34 41 48

- Create the following R vector using **only** the **seq()** and **rep()** functions:
  -8 -8 -8 -8 0 8 8 8 16 16 16 16 16

## 3.2 Operations on vectors: element-wise

- All operations on vectors in R happen **element by element** (cfr. *NumPy*).

- **Vector Recycling**:

  If 2 vectors of **different** lengths are involved in an operation, the **shortest vector** will be repeated until all elements of the longest vector are matched.
  A *warning* message will be sent to the stdout.

### 3.2.1 Examples

```r
x <- -3:3
x
```

```
[1] -3 -2 -1  0  1  2  3
```

```r
y <- 1:7
y
```

```
[1] 1 2 3 4 5 6 7
```

```
xy <- x*y
xy
```

```
[1] -3 -4 -3  0  5 12 21
```

```
xpy <- x^y
xpy
```

```
[1]   -3    4   -1    0    1   64 2187
```

```
x <- 0:10
y <- 1:2
length(x)
```

```
[1] 11
```

```
length(y)
```

```
[1] 2
```

```
x
```

```
 [1]  0  1  2  3  4  5  6  7  8  9 10
```

```
y
```

```
[1] 1 2
```

```
x+y
```

```
Warning in x + y: longer object length is not a multiple of shorter object
length
```

```
 [1]  1  3  3  5  5  7  7  9  9 11 11
```

### 3.2.2 Exercises

- Create the following vector (do **not** use **c()**!):
  -512 -216 -64 -8 0 8 64 216 512 1000

## 3.3 Retrieving elements of vectors

- Indexing: starts at **1** (**not 0** like C/C++, Python, Java, . . . .) see also:
  Edsger Dijkstra: Why numbering should start at zero
- Use of vector with indices to extract values.
- Advanced features:
  - use of boolean values to extract values.
  - the membership operator: **%in%**.
  - the deselect/omit operator: **-**
  - **which()**: returns the indices for which the condition is true.
  - **any()**/**all()** functions.
    * **any()** : **TRUE** if at least 1 value is true
    * **all()** : **TRUE** if all values are true

### 3.3.1 Examples

- Use of a simple index:

```
x <- seq(2,100,by=15)
x
```

```
[1]   2 17 32 47 62 77 92
```

```
x[4]
```

```
[1] 47
```

```
x[1]
```

```
[1] 2
```

- Select several indices at once using vectors:

```
x
```

```
[1]   2 17 32 47 62 77 92
```

```
x[3:5]
```

```
[1] 32 47 62
```

```
x[c(1,3,5,7)]
```

```
[1]   2 32 62 92
```

```
x[seq(1,7,by=2)]
```

```
[1]   2 32 62 92
```

- Extraction via booleans (i.e. retain only those values that are equal to **TRUE**):

```
x
```

```
[1]   2 17 32 47 62 77 92
```

```
x>45
```

```
[1] FALSE FALSE FALSE  TRUE  TRUE  TRUE  TRUE
```

```
x[x>45]
```

```
[1] 47 62 77 92
```

- Use of the **%in%** operator:

```
x
```

```
[1]   2 17 32 47 62 77 92
```

```
10 %in% x
```

```
[1] FALSE
```

```r
62 %in% x
```

```
[1] TRUE
```

```r
c(32,33,43) %in% x
```

```
[1]  TRUE FALSE FALSE
```

```r
!(c(32,33,43) %in% x)
```

```
[1] FALSE  TRUE  TRUE
```

- Negate/filter out the elements with **negative** indices:

```r
x <- c(1,13,17,27,49,91)
x
```

```
[1]  1 13 17 27 49 91
```

```r
x[-c(2,4,6)]
```

```
[1]  1 17 49
```

```r
z <- x[-1] - x[-length(x)]
z
```

```
[1] 12  4 10 22 42
```

- The **which()** function returns **only those indices** of which the condition/expression is **true**.

```r
# Sample 10 numbers from N(0,1)
vecnum <- rnorm(n=10)
vecnum
```

```
 [1]  1.21026318 -0.97310436 -0.01271416 -1.37277634 -0.13987563  0.41705537
 [7] -1.02959734 -0.86444438 -0.71286292  0.52680419
```

```r
which(vecnum>1.0)
```

```
[1] 1
```

- Use of the **any()**/**all()** functions.

```r
y <- seq(0,100,by=10)
x
```

```
[1]  1 13 17 27 49 91
```

```r
y
```

```
 [1]   0  10  20  30  40  50  60  70  80  90 100
```

```r
any(x<y)
```

```
Warning in x < y: longer object length is not a multiple of shorter object
length
```

```
[1] TRUE
```
```
all(x[6:7]>y[2:3])
```
```
[1] NA
```

### 3.3.2 Exercises

- R has the its own inversion function, **rev()**, e.g.:

```
x <- seq(from=2,to=33,by=3)
x
```
```
 [1]  2  5  8 11 14 17 20 23 26 29 32
```
```
y <- rev(x)
y
```
```
 [1] 32 29 26 23 20 17 14 11  8  5  2
```

  Invert the vector x without invoking the **rev()** function.

- The `Taylor series` for $\ln(1+x)$ is converging when $|x| < 1$ and is given by:

$$\ln(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \frac{x^5}{5} - \frac{x^6}{6} + \dots$$

  Calculate the sum of the first 5, 10, 15 terms in the above expression to approximate $\ln(1.2)$. Compare with R's value i.e.: **log(1.2)**.

- The `logarithmic` return in finance is defined as:

$$R_t = \ln\left(\frac{P_t}{P_{t-1}}\right)$$

  – Generate a financial time series using the following `R` code:

```
price <- abs(rcauchy(1000))+1.E-6
```

  – Calculate the `logarithmic return` for the financial time series `price`.
    The newly created time series will be 1 element shorter in length than the original one.
    Compare your result with **diff(log(price))**.

- `Monte-Carlo` approximation of $\pi$

  Let `S1` be the square spanned by the following 4 vertices: $\{(0,0),(0,1),(1,0),(1,1)\}$.
  Let `S2` be the first quadrant of the unit-circle $\mathcal{C} : x^2 + y^2 = 1$.

  The ratio $\rho$ defined as:
$$\rho := \frac{\text{Area S2}}{\text{Area S1}} = \frac{\#\text{Points in S2}}{\#\text{Points in S1}}$$

  allows us to estimate $\frac{\pi}{4}$ numerically.

  Therefore:

  – Sample 100000 independent $x$-coordinates from `Unif`.
  – Sample 100000 independent $y$-coordinates from `Unif`.
  – Calculate an approximate value for $\pi$ using the Monte-Carlo approach.

  Note: The uniform distribution $[0,1)$ (`Unif`) can be sampled using **runif()**.

## 3.4 Hash tables

A **hash table** is a data structure which implements an associative array or dictionary. It is an abstract data which maps data to keys.

- There are several ways to create one:
  - Map names to an existing vector
  - Add names when creating the vector
- To remove the map, map the names to NULL

### 3.4.1 Examples

- Creation of 2 independent vectors

```
capitals <- c("Albany", "Providence", "Hartford", "Boston", "Montpelier", "Concord", "Augusta")
states <- c("NY", "RI", "CT", "MA", "VT", "NH", "ME")
capitals
```

```
[1] "Albany"     "Providence" "Hartford"   "Boston"     "Montpelier"
[6] "Concord"    "Augusta"
```

```
states
```

```
[1] "NY" "RI" "CT" "MA" "VT" "NH" "ME"
```

```
capitals[3]
```

```
[1] "Hartford"
```

- Create the hashtable/dictionary

```
# Method 1
names(capitals) <- states
capitals
```

```
         NY           RI           CT           MA           VT           NH
   "Albany" "Providence"   "Hartford"     "Boston" "Montpelier"    "Concord"
         ME
  "Augusta"
```

```
capitals["MA"]
```

```
      MA
"Boston"
```

```
names(capitals)
```

```
[1] "NY" "RI" "CT" "MA" "VT" "NH" "ME"
```

```
# Method 2
phonecode <- c("801"="SLC", "206"="Seattle", "307"="Wyoming")
phonecode
```

```
      801       206       307
    "SLC" "Seattle" "Wyoming"
```

```
phonecode["801"]
```

```
  801
"SLC"
```

- Dissociate the 2 vectors

```
names(capitals) <- NULL
capitals
```

```
[1] "Albany"     "Providence" "Hartford"   "Boston"      "Montpelier"
[6] "Concord"    "Augusta"
```

## 3.5  NA (Not Available values)

- **NA**: stands for 'Not Available'/Missing values and has a length of 1.
  There are in essence 4 versions depending on the type:

  – **NA** (logical - **default**)
  – **NA_integer** (integer)
  – **NA_real** (double precision)
  – **NA_character** (string)

  Under the hood, the version of NA is subjectd to **coercion**:
  *logical → integer → double → character*

- some functions e.g. **mean()** return (by default) NA if
  1 or more instances NA are present in a vector.

- **is.na()**: test a vector (element-wise) for NA values.
  **Do NOT use:**

```
x == NA
```

  **but use INSTEAD:**

```
is.na(x)
```

### 3.5.1  Examples

- Types of NA

```
x <- NA
typeof(x)
```

```
[1] "logical"
```

```
# logical NA coerced to double precision NA
x <- c(3.0, 5.0, NA)
typeof(x[3])
```

```
[1] "double"
```

\* Functions on a vector containing NA

```
mean(x)
```

```
[1] NA
```
```
mean(x, na.rm=TRUE)
```

```
[1] 4
```

\* Check of the NA availability

```
x <- c(NA, 1, 2, NA)
is.na(x)
```

```
[1]  TRUE FALSE FALSE  TRUE
```

\* Functions on a vector containing NA

```
mean(x)
```

```
[1] NA
```
```
mean(x, na.rm=TRUE)
```

```
[1] 1.5
```

### 3.5.2   Exercises

- A family has installed a device to monitor their daily energy consumption (in kWh).
  When a measurement fails or is unavailable NA is recorded.

  You can invoke the following code to generate the measurements generated by the device.

  ```
  dailyusage <- 30.0 + runif(365, min=0, max=5.0)
  dailyusage[sample(1:365, sample(1:50,1), replace=FALSE)] <- NA
  ```

  - How many measurements failed?
  - What is the average daily energy consumption (based on the non-failed) measurements?

## 3.6   NaN and infinities

- **NaN** (only for numeric types!), and the infinties **Inf** and **-Inf**
  are part of the IEEE 754 floating-point standard.
- To test whether you have:
  - finite numbers: use **is.finite()**
  - infinite numbers: use **is.infinite()**
  - NaNs: use **is.nan()**
- Further:
  - a **NaN** will return **TRUE** when tested by either **is.nan()** or **is.na()**
  - a **NA** will return **TRUE** only when tested by **is.na()**

### 3.6.1   Examples

- Infinities:

```r
x <- 5.0/0.0
x
```

```
[1] Inf
```

```r
is.finite(x)
```

```
[1] FALSE
```

```r
is.infinite(x)
```

```
[1] TRUE
```

```r
is.nan(x)
```

```
[1] FALSE
```

```r
y <- -5.0/0.0
y
```

```
[1] -Inf
```

```r
is.finite(y)
```

```
[1] FALSE
```

```r
is.infinite(y)
```

```
[1] TRUE
```

```r
is.nan(y)
```

```
[1] FALSE
```

```r
z <- x + y
z
```

```
[1] NaN
```

```r
typeof(z)
```

```
[1] "double"
```

```r
is.finite(z)
```

```
[1] FALSE
```

```r
is.infinite(z)
```

```
[1] FALSE
```

```r
is.nan(z)
```

```
[1] TRUE
```

- **is.na()** vs. **is.nan()**:

```
# is.nan
v <- c(NA, z, 5.0, log(-1.0))
```

```
Warning in log(-1): NaNs produced
is.nan(v)
```

```
[1] FALSE  TRUE FALSE  TRUE
```

```
# is.na(): also includes NaN!
v <- c(NA, z, 5.0, log(-1.0))
```

```
Warning in log(-1): NaNs produced
is.na(v)
```

```
[1]  TRUE  TRUE FALSE  TRUE
```

## 3.7   Note on logical operators

- **&**, **|**, **!**, **xor()**: **element-wise** operators on vectors (cfr. arithmetic operators)
- **&&**, **||**: evaluated from **left** to **right** until result is determined.

### 3.7.1   Examples

- Vector operators (**&**, **|**, **!** and **xor()**)

```
x <- sample(x=1:10, size=10, replace=TRUE)
x
```

```
 [1]  3  3  4 10  2  3  5  4  6  4
```

```
y <- sample(x=1:10, size=10, replace=TRUE)
y
```

```
 [1] 4 5 8 1 6 7 6 2 2 2
```

```
v1 <- (x<=3)
v1
```

```
 [1]  TRUE  TRUE FALSE FALSE  TRUE  TRUE FALSE FALSE FALSE FALSE
```

```
v2 <- (y>=7)
v2
```

```
 [1] FALSE FALSE  TRUE FALSE FALSE  TRUE FALSE FALSE FALSE FALSE
```

```
v1 & v2
```

```
[1] FALSE FALSE FALSE FALSE FALSE  TRUE FALSE FALSE FALSE FALSE
```

```
v1 | v2
```

```
[1]  TRUE  TRUE  TRUE FALSE  TRUE  TRUE FALSE FALSE FALSE FALSE
```

```
xor(v1, v2)
```

```
[1]  TRUE  TRUE  TRUE FALSE  TRUE FALSE FALSE FALSE FALSE FALSE
```

```
!v1
```

```
[1] FALSE FALSE  TRUE  TRUE FALSE FALSE  TRUE  TRUE  TRUE  TRUE
```

### 3.7.2   Exercises

- Generate a random vector of integers using the following code:

  ```
  x <- sample(x=0:1000,size=100, replace=TRUE)
  ```

  - Invoke the above code to generate the vector x
  - Find if there are any integers in the vector x which can be divided by 4 and 6
  - Find those numbers and their corresponding indices in the vector x.

```
"It is my experience that proofs involving matrices can be
shortened by 50% if one throws the matrices out" (Emil Artin)
```

# 4   Matrices & Arrays

Matrices and arrays are **homogeneous atomic vectors** with
an **extra** attribute: dimension

By default, the elements are stored in a **column-major** fashion. (cfr. **Fortran**).
However, we can store the elements in **row-major** order (cfr. **C**) as well.

## 4.1   Creation of matrices

Matrices can be created in different ways:

- use of the **matrix()** function
- use of **rbind()**/**cbind()**
- set the **attributes()** of a vector
- special functions like e.g. **diag()**

### 4.1.1   Examples

- use of the **matrix()** function:

The **matrix()** function creates a matrix based on a vector.
By default, the elements are stored in a **column-major** fashion.

The use of the flag **byrow=TRUE** will store the data in a **row-major** fashion.

```
A <- matrix(data=1:10, nrow=2)    # Column-major (like Fortran)
A

     [,1] [,2] [,3] [,4] [,5]
[1,]    1    3    5    7    9
[2,]    2    4    6    8   10
B <- matrix(data=c(2,3,893,0.17), nrow=2, ncol=2)
B

     [,1]   [,2]
[1,]    2 893.00
[2,]    3   0.17



C <- matrix(data=1:10, nrow=2, byrow=TRUE)    # Row-major (like C, C++)
C

     [,1] [,2] [,3] [,4] [,5]
[1,]    1    2    3    4    5
[2,]    6    7    8    9   10
```

- use of the **rbind()**/**cbind()** functions:

- **rbind()**: Bind several vectors (as rows) into a matrix.
- **cbind()**: Bind several vectors (as columns) into a matrix.

```
A <- rbind(1:10,11:20)
A
```

```
     [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
[1,]    1    2    3    4    5    6    7    8    9    10
[2,]   11   12   13   14   15   16   17   18   19    20
```

```
typeof(A)
```

```
[1] "integer"
```

```
class(A)
```

```
[1] "matrix" "array"
```

```
B <- cbind(1:5,6:10,11:15)
B
```

```
     [,1] [,2] [,3]
[1,]    1    6   11
[2,]    2    7   12
[3,]    3    8   13
[4,]    4    9   14
[5,]    5   10   15
```

```
class(B)
```

```
[1] "matrix" "array"
```

### 4.1.2 Matrices: vectors with a non-NULL dim attribute

The **fundamental** difference between an `R` vector and matrix is the presence (in the case of matrices) of a non `NULL` **dim** attribute.

We can easily convert a vector into a matrix by setting the dimensions of the vector:

- through the **dim()** function.
- through the **attr()** function.

The inverse can be done as well by setting the **dim** attribute of matrix to `NULL`.

```
A <- 1:10
typeof(A)
```

```
[1] "integer"
```

```
class(A)
```

```
[1] "integer"
```

```
dim(A)
```

```
NULL
```

```
# Matrix
B <- matrix(1:10,nrow=2,ncol=5,byrow=TRUE)
typeof(B)
```

```
[1] "integer"
```

```
class(B)
```

```
[1] "matrix" "array"
```

```
dim(B)
```

```
[1] 2 5
```

```
# Vector
A <- 1:10
A
```

```
 [1]  1  2  3  4  5  6  7  8  9 10
```

```
dim(A)
```

```
NULL
```

```
typeof(A)
```

```
[1] "integer"
```

```
class(A)
```

```
[1] "integer"
```

```
# OPTION I: Using the dim function transform a vector into a matrix
dim(A) <- c(2,5)
A
```

```
     [,1] [,2] [,3] [,4] [,5]
[1,]    1    3    5    7    9
[2,]    2    4    6    8   10
```

```
dim(A)
```

```
[1] 2 5
```

```
typeof(A)
```

```
[1] "integer"
```

```
class(A)
```

```
[1] "matrix" "array"
```

```
# Converting the matrix back to a vector
dim(A) <- NULL
dim(A)
```

NULL

```
typeof(A)
```

[1] "integer"

```
class(A)
```

[1] "integer"

```
# Option II: More general way
# Convert vector into a matrix
A <- 1:8
A
```

[1] 1 2 3 4 5 6 7 8

```
class(A)
```

[1] "integer"

```
attr(A,'dim') <- c(2,4)
A
```

```
     [,1] [,2] [,3] [,4]
[1,]    1    3    5    7
[2,]    2    4    6    8
```

```
class(A)
```

[1] "matrix" "array"

```
# Convert matrix into a vector.
attr(A, 'dim') <- NULL
A
```

[1] 1 2 3 4 5 6 7 8

```
class(A)
```

[1] "integer"

## 4.2 Retrieving elements/subsetting

Matrices (and arrays) can be subsetted in different ways:

- use an **index** for each dimension, where the dimensions are comma-separated
  - If an **index** for a dimension is **omitted**:
    consider all dimensions (may lead to reduction of the dimension)
  - **but** you can use **drop=FALSE** to prevent dimensionality reduction.
- use another **vector** (can be either linear or a vector for each dimension)

- by using another **matrix**.

### 4.2.1 Examples

- Use of **indices**:

```
A <- matrix(1:30, nrow=6, ncol=5)
A
```

```
     [,1] [,2] [,3] [,4] [,5]
[1,]    1    7   13   19   25
[2,]    2    8   14   20   26
[3,]    3    9   15   21   27
[4,]    4   10   16   22   28
[5,]    5   11   17   23   29
[6,]    6   12   18   24   30
```

```
A[3,4]
```

```
[1] 21
```

```
A[6,2]
```

```
[1] 12
```

```
x1 <- A[2,]
x1
```

```
[1]  2  8 14 20 26
```

```
dim(x1)
```

```
NULL
```

```
x2 <- A[,3]
x2
```

```
[1] 13 14 15 16 17 18
```

```
dim(x2)
```

```
NULL
```

The flag **drop=FALSE** can be used to prevent dimensionality reduction

```
y1 <- A[2,,drop=FALSE]
y1
```

```
     [,1] [,2] [,3] [,4] [,5]
[1,]    2    8   14   20   26
```

28

```
dim(y1)
```

```
[1] 1 5
```

```
y2 <- A[,3,drop=FALSE]
y2
```

```
     [,1]
[1,]   13
[2,]   14
[3,]   15
[4,]   16
[5,]   17
[6,]   18
```

```
dim(y2)
```

```
[1] 6 1
```

- Use of **vector(s)**:

```
A
```

```
     [,1] [,2] [,3] [,4] [,5]
[1,]    1    7   13   19   25
[2,]    2    8   14   20   26
[3,]    3    9   15   21   27
[4,]    4   10   16   22   28
[5,]    5   11   17   23   29
[6,]    6   12   18   24   30
```

```
x1 <- A[2:4,]
x1
```

```
     [,1] [,2] [,3] [,4] [,5]
[1,]    2    8   14   20   26
[2,]    3    9   15   21   27
[3,]    4   10   16   22   28
```

```
dim(x1)
```

```
[1] 3 5
```

```
x2 <- A[,1:3]
x2
```

```
     [,1] [,2] [,3]
[1,]    1    7   13
[2,]    2    8   14
[3,]    3    9   15
[4,]    4   10   16
[5,]    5   11   17
```

```
[6,]     6    12    18
dim(x2)

[1] 6 3
```

```
# Using a vector for EACH dimension
A[c(1,3),c(2,4)]

     [,1] [,2]
[1,]    7   19
[2,]    9   21
```

```
# Using 1 vector => Linear index
A[c(1,3,8,10)]

[1]  1  3  8 10
```

```
A[c(TRUE,FALSE,TRUE,TRUE,FALSE,TRUE),c(2,3)]

     [,1] [,2]
[1,]    7   13
[2,]    9   15
[3,]   10   16
[4,]   12   18
```

```
A[c(TRUE,FALSE,TRUE,TRUE,FALSE,TRUE),]

     [,1] [,2] [,3] [,4] [,5]
[1,]    1    7   13   19   25
[2,]    3    9   15   21   27
[3,]    4   10   16   22   28
[4,]    6   12   18   24   30
```

```
# Use of a linear index
A[c(TRUE,FALSE,TRUE,TRUE,FALSE,TRUE)]

 [1]  1  3  4  6  7  9 10 12 13 15 16 18 19 21 22 24 25 27 28 30
```

- Use of a **matrix**:

```
A

     [,1] [,2] [,3] [,4] [,5]
[1,]    1    7   13   19   25
[2,]    2    8   14   20   26
[3,]    3    9   15   21   27
[4,]    4   10   16   22   28
```

```
[5,]    5   11   17   23   29
[6,]    6   12   18   24   30
```

```
mysubset <- matrix(c( 2, 1,
                      3, 5,
                      4, 2,
                      6, 5), ncol=2, byrow=TRUE )
A[mysubset]
```

```
[1]  2 27 10 30
```

### 4.2.2 Exercises

- Create the following matrix A, given by:

```
      [,1] [,2] [,3]  [,4]     [,5]      [,6]
[1,]    3    9   27    81      243       729
[2,]    5   25  125   625     3125     15625
[3,]    7   49  343  2401    16807    117649
[4,]   11  121 1331 14641   161051   1771561
[5,]   13  169 2197 28561   371293   4826809
[6,]   17  289 4913 83521  1419857  24137569
```

  1. get element 343

  2. get the elements 25, 625, 2197 and 4826809 (all at once).

  3. get the fourth row as a vector.

  4. get the fourth row as a matrix.

  5. get columns 2 and 3 (at the same time).

  6. get everything except rows 2 and 4.

  7. the diagonal of matrix A.

## 4.3 Operations on matrices

- Operations like **\*,/, +** happpen element-wise.
- There are also more specialized functions:
  - the mean over rows and columns (**rowMeans()**, **colMeans()**)
  - linear algebra functions (**%\*%**, **t()**, ...)

### 4.3.1 Examples

- Operations (by **default: element-by-element**):

```
A <- matrix(1:10, nrow=2)
B <- matrix( seq(10, 100, by=10), nrow=2)
A
```

```
     [,1] [,2] [,3] [,4] [,5]
[1,]    1    3    5    7    9
[2,]    2    4    6    8   10
```

31

```
B
```

```
     [,1] [,2] [,3] [,4] [,5]
[1,]   10   30   50   70   90
[2,]   20   40   60   80  100
```

```
A*B
```

```
     [,1] [,2] [,3] [,4] [,5]
[1,]   10   90  250  490  810
[2,]   40  160  360  640 1000
```

```
C <- matrix(rep(2,10), nrow=2)
C
```

```
     [,1] [,2] [,3] [,4] [,5]
[1,]    2    2    2    2    2
[2,]    2    2    2    2    2
```
```
C**A
```

```
     [,1] [,2] [,3] [,4] [,5]
[1,]    2    8   32  128  512
[2,]    4   16   64  256 1024
```

- Calculate **row and column means** :

```
# Means of rows and columns
A
```

```
     [,1] [,2] [,3] [,4] [,5]
[1,]    1    3    5    7    9
[2,]    2    4    6    8   10
```

```
rowMeans(A)
```

```
[1] 5 6
```

```
colMeans(A)
```

```
[1] 1.5 3.5 5.5 7.5 9.5
```

- **Matrix multiplication** (%*%) :

```
A <- matrix(1:6, nrow=2)
A
```

```
     [,1] [,2] [,3]
[1,]    1    3    5
[2,]    2    4    6
```

```
B <- matrix(seq(10,120,by=10), nrow=3)
B
```

```
     [,1] [,2] [,3] [,4]
[1,]   10   40   70  100
[2,]   20   50   80  110
```

```
[3,]   30   60   90  120
```

```
C <- A%*%B
C
```

```
      [,1] [,2] [,3] [,4]
[1,]   220  490  760 1030
[2,]   280  640 1000 1360
```

```
dim(C)
```

```
[1] 2 4
```

- **Linear algebra** routines

Some of the more common ones in R:

- **solve()** : **invert** a square matrix
- **diag()**
  - **extracts** the diagonal of a matrix when a matrix is provided.
  - **creates** a diagonal matrix when a vector is provided.
- **eigen()** : calculates the **eigenvalues** and **eigenvectors** of a matrix
- **det()** : calculates the **determinant** of a matrix.
- **t()**: calculates the **transpose**[1] of a matrix.

```
# Invert matrix A
A <- matrix(c(1, 3, 2, 4), ncol = 2, byrow = T)
Ainv <- solve(A)
Ainv %*% A
```

```
     [,1] [,2]
[1,]    1    0
[2,]    0    1
```

```
# Create a diagonal matrix
C <- diag(c(1,4,7))
C
```

```
     [,1] [,2] [,3]
[1,]    1    0    0
[2,]    0    4    0
[3,]    0    0    7
```

```
# Extract the diagonal elements
D <- matrix(1:8,nrow=4)
D
```

```
     [,1] [,2]
[1,]    1    5
[2,]    2    6
[3,]    3    7
[4,]    4    8
```

---

[1]Can also be used for dataframes (see later)

```
diag(D)
```

```
[1] 1 6
```

```
# Calculate eigenvalues and eigenvectors of A
r <- eigen(A)
r
```

```
eigen() decomposition
$values
[1]  5.3722813 -0.3722813

$vectors
            [,1]       [,2]
[1,] -0.5657675 -0.9093767
[2,] -0.8245648  0.4159736
```

```
# Eigenvalues
r$values
```

```
[1]  5.3722813 -0.3722813
```

```
# Matrix with eigenvectors
r$vectors
```

```
            [,1]       [,2]
[1,] -0.5657675 -0.9093767
[2,] -0.8245648  0.4159736
```

```
# Diagonal Matrix (Similarity Transformation)
solve(r$vectors) %*% A %*% r$vectors
```

```
              [,1]        [,2]
[1,]  5.372281e+00  0.0000000
[2,] -3.330669e-16 -0.3722813
```

Note that under the hood R calls BLAS and LAPACK.

```
# Find the version used of BLAS and LAPACK
La_library()
```

```
[1] "/usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.7.1"
```

```
extSoftVersion()["BLAS"]
```

```
                                                    BLAS
"/usr/lib/x86_64-linux-gnu/blas/libblas.so.3.7.1"
```

### 4.3.2 Exercises

- Linear regression:

    - **Step 1**:
      Create a **synthetic** data set by executing the following R code:

```
x <- seq(from=0, to=20.0, by=0.25)
a <- 2.0
b <- 1.5
c <- 0.5
y <- a + b*x + c*x^2 + rnorm(length(x))
```

– **Step 2**:
Our goal is to use the following linear model, i.e.:

$$Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \epsilon_i$$

or in matrix form:

$$Y = X\beta + \epsilon \tag{1}$$

to fit the previously generated data set.

In Eq.(1), we have:

* $Y$ : a $n \times 1$ column vector.

* $X$ : a $n \times 3$ matrix.

* $\beta$: a $3 \times 1$ column vector.

* $\epsilon$ is : a $n \times 1$ column vector and $\sim N(0, \sigma^2)$

An estimate for $\beta$ $(\widehat{\beta})$ can be found (using Least-Squares, MLE see e.g. (Seber & Lee, 2012)) and has the following form:

$$\widehat{\beta} = (\mathbf{X^T X})^{-1}\mathbf{X^T Y} \tag{2}$$

where:

the column vector $\mathbf{Y}$ is given by:

$$\mathbf{Y} := \begin{bmatrix} y_1 \\ y_2 \\ .. \\ y_n \end{bmatrix}$$

and the matrix $\mathbf{X}^2$ takes the following form:

$$\mathbf{X} := \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{bmatrix}$$

**Calculate** $\widehat{\beta}$ using Eq.(2).

An estimate for the residuals $(\widehat{\epsilon})$ is given by:

$$\widehat{\epsilon} = \mathbf{Y} - \mathbf{X}\widehat{\beta} \tag{3}$$

**Calculate** $\widehat{\epsilon}$ using Eq.(3).

– **Step 3**:

You can check your results using the following R code.

---
[2]This is a known as a Vandermonde matrix.

35

```
                 myquadfit <- lm(y ~ x + I(x^2))
                 cat(sprintf("The estimates for beta::\n"))
                 cat(myquadfit$coefficients)
                 cat(sprintf("The residuals::\n"))
                 cat(myquadfit$residuals)
```

## 4.4   Hash tables/dictionaries

We can also use hashes for matrices. We can select one or both dimensions. To create hashes, for: - rows: use **rownames** - columns: use **colnames**

To remove the hash, use the **NULL** (like for vectors).

### 4.4.1   Examples

```
A1 <- c(0      , 5471.52, 5091.57, 5392.82,
        5416.45, 4584.33, 4904.83, 3851.73)
A2 <- c(5471.52,       0, 1315.28,  927.35,
        1505.11,  944.40, 1157.42, 1945.42)
A3 <- c(5091.57, 1315.28,       0, 2166.00,
        2724.01, 1571.76,  293.52, 1240.77)
A4 <- c(5392.82,  927.35, 2166.00,       0,
         577.85,  973.23, 1947.28, 2422.32)
A5 <- c(5416.45, 1505.11, 2724.01,  577.85,
              0, 1366.63, 2490.97, 2838.62)
A6 <- c(4584.33,  944.40, 1571.76,  973.23,
        1366.63,       0, 1290.15, 1474.26)
A7 <- c(4904.83, 1157.42,  293.52, 1947.28,
        2490.97, 1290.15,       0, 1064.41)
A8 <- c(3851.73, 1945.42, 1240.77, 2422.32,
        2838.62, 1474.26, 1064.41,       0)
```

```
dist <- rbind(A1,A2,A3,A4,A5,A6,A7,A8)
dist

      [,1]    [,2]    [,3]    [,4]    [,5]    [,6]    [,7]    [,8]
A1    0.00 5471.52 5091.57 5392.82 5416.45 4584.33 4904.83 3851.73
A2 5471.52    0.00 1315.28  927.35 1505.11  944.40 1157.42 1945.42
A3 5091.57 1315.28    0.00 2166.00 2724.01 1571.76  293.52 1240.77
A4 5392.82  927.35 2166.00    0.00  577.85  973.23 1947.28 2422.32
A5 5416.45 1505.11 2724.01  577.85    0.00 1366.63 2490.97 2838.62
A6 4584.33  944.40 1571.76  973.23 1366.63    0.00 1290.15 1474.26
A7 4904.83 1157.42  293.52 1947.28 2490.97 1290.15    0.00 1064.41
A8 3851.73 1945.42 1240.77 2422.32 2838.62 1474.26 1064.41    0.00
```

```
# Adding hashes to both rows and columns
cities <- c("Anchorage","Atlanta","Austin","Baltimore","Boston", "Chicago", "Dallas","Denver")
rownames(dist) <- cities
```

```
colnames(dist) <- cities
dist
```

```
          Anchorage Atlanta  Austin Baltimore  Boston Chicago  Dallas  Denver
Anchorage      0.00 5471.52 5091.57   5392.82 5416.45 4584.33 4904.83 3851.73
Atlanta     5471.52    0.00 1315.28    927.35 1505.11  944.40 1157.42 1945.42
Austin      5091.57 1315.28    0.00   2166.00 2724.01 1571.76  293.52 1240.77
Baltimore   5392.82  927.35 2166.00      0.00  577.85  973.23 1947.28 2422.32
Boston      5416.45 1505.11 2724.01    577.85    0.00 1366.63 2490.97 2838.62
Chicago     4584.33  944.40 1571.76    973.23 1366.63    0.00 1290.15 1474.26
Dallas      4904.83 1157.42  293.52   1947.28 2490.97 1290.15    0.00 1064.41
Denver      3851.73 1945.42 1240.77   2422.32 2838.62 1474.26 1064.41    0.00
```

```
dist["Chicago", "Denver"]
```

```
[1] 1474.26
```

```
dist["Austin", "Boston"]
```

```
[1] 2724.01
```

## 4.5 Arrays

Say something about arrays.

# 5 Special Data Types (Factors and Date/Time types)

Every R object has attributes (i.e. properties or metadata).
They can be classified as:

- **intrinisic** properties e.g. **length()**
- **external** properties (to be set by the user)

## 5.1 Attributes

- can be get/retrieved using **attributes()**.
- can be set:
    - individually using **attr()**
    - in generally using **structure()**
- some attributes can (also) be set/unset with **special** functions:
    - names: **names()**
    - dimension: **dim()**
    - comment : **comment()**
    - time series: **tsp()**
    - factor : **factor()** (see next section)

### 5.1.1 Examples

- 1 attribute:

```r
x <- 1:5
x
```

```
[1] 1 2 3 4 5
```

```r
attr(x, 'prop1') <- "hello"
attributes(x)
```

```
$prop1
[1] "hello"
```

```r
x
```

```
[1] 1 2 3 4 5
attr(,"prop1")
[1] "hello"
```

```r
attr(x, 'prop1') <- NULL
attributes(x)
```

```
NULL
```

```r
x
```

```
[1] 1 2 3 4 5
```

- more than 1 attribute:

```r
y <- 1:8
y
```

```
[1] 1 2 3 4 5 6 7 8
```

```r
y <- structure(y, dim=c(2,4), tag="trial")
y
```

```
     [,1] [,2] [,3] [,4]
[1,]    1    3    5    7
[2,]    2    4    6    8
attr(,"tag")
[1] "trial"
```

```r
attributes(y)
```

```
$dim
[1] 2 4

$tag
[1] "trial"
```

```r
typeof(y)
```

```
[1] "integer"
```

```r
class(y)
```

```
[1] "matrix" "array"
```

```r
# Remove BOTH attributes
y <- structure(y, dim=NULL, tag=NULL)
y
```

```
[1] 1 2 3 4 5 6 7 8
```

```r
attributes(y)
```

```
NULL
```

```r
typeof(y)
```

```
[1] "integer"
```

```r
class(y)
```

```
[1] "integer"
```

- **names()**

```r
# Set the names attribute
capitals <- c("Salt Lake City", "Carson City", "Boise", "Santa Fe")
```

```r
names(capitals) <- c("UT", "NV", "ID", "NM")
capitals
```

```
            UT               NV               ID               NM
"Salt Lake City"    "Carson City"          "Boise"       "Santa Fe"
```

```r
attributes(capitals)
```

```
$names
[1] "UT" "NV" "ID" "NM"
```

```r
# Remove the names attribute
names(capitals) <- NULL
capitals
```

```
[1] "Salt Lake City" "Carson City"    "Boise"          "Santa Fe"
```

- **dim()**

```r
x <- 1:12
x
```

```
 [1]  1  2  3  4  5  6  7  8  9 10 11 12
```

```r
typeof(x)
```

```
[1] "integer"
```

```r
class(x)
```

```
[1] "integer"
```

```r
# Set the dimension attribute
dim(x) <- c(3,4)
x
```

```
     [,1] [,2] [,3] [,4]
[1,]    1    4    7   10
[2,]    2    5    8   11
[3,]    3    6    9   12
```

```r
typeof(x)
```

```
[1] "integer"
```

```r
class(x)
```

```
[1] "matrix" "array"
```

```
# Remove the dimension attribute
dim(x) <- NULL
x
```

```
 [1]  1  2  3  4  5  6  7  8  9 10 11 12
```
```
typeof(x)
```

```
[1] "integer"
```
```
class(x)
```

```
[1] "integer"
```

- **comment()**
```
x <- structure(1:6, comment="My vector")
typeof(x)
```

```
[1] "integer"
```
```
class(x)
```

```
[1] "integer"
```
```
comment(x)
```

```
[1] "My vector"
```

## 5.2 Factor variables (Categorical variables)

- Factor variables (factors, categorical variables) are discrete variables (i.e not continuous).
  The factors bear labels (**levels**) which are mapped into **integers**.

- Therefore, factors are stored as integer vector with 2 attributes:

    - **class**= "factor"
    - **levels**: a vector with the "labels".

- By default (**unordered**) the labels are mapped **alphabetically** to the integers. We can **impose** our own **ordering** between integers and labels (levels).

- Useful functions:

    - **levels()** : provides the levels of a factor
    - **table()**: returns the counts of each level
    - **is.factor()**: tests whether a variable is a factor variable
    - **is.ordered()**: tests whether a variable is an ordered factor variable

### 5.2.1 Examples

- Creation of an **unordered** factor
```
# Creation of an unordered factor
temp.data <- c("High","Low","VeryHigh","Low","VeryLow","Medium",
               "VeryHigh","VeryHigh","Low","Low","Medium","VeryHigh",
               "VeryHigh","VeryHigh","Low","High","VeryLow")
myfac.temp.data <- factor(temp.data)
myfac.temp.data
```

41

```
[1] High      Low       VeryHigh Low       VeryLow  Medium    VeryHigh VeryHigh
[9] Low       Low       Medium   VeryHigh VeryHigh VeryHigh Low       High
[17] VeryLow
Levels: High Low Medium VeryHigh VeryLow
```

```r
# by default: the levels are stored ALPHABETICALLY (i.e. unordered)
levels(myfac.temp.data)
```

```
[1] "High"     "Low"      "Medium"   "VeryHigh" "VeryLow"
```

```r
table(myfac.temp.data)
```

```
myfac.temp.data
    High      Low   Medium VeryHigh  VeryLow
       2        5        2        6        2
```

```r
is.factor(myfac.temp.data)
```

```
[1] TRUE
```

```r
is.ordered(myfac.temp.data)
```

```
[1] FALSE
```

- Creation of an **ordered** factor

```r
# Creation of an unordered factor
temp.data <- c("High","Low","VeryHigh","Low","VeryLow","Medium",
               "VeryHigh","VeryHigh","Low","Low","Medium","VeryHigh",
               "VeryHigh","VeryHigh","Low","High","VeryLow")
myfac2.temp.data <- factor(temp.data, ordered=TRUE,
                           levels=c("VeryLow","Low","Medium","High","VeryHigh"))
myfac2.temp.data
```

```
[1] High      Low       VeryHigh Low       VeryLow  Medium    VeryHigh VeryHigh
[9] Low       Low       Medium   VeryHigh VeryHigh VeryHigh Low       High
[17] VeryLow
Levels: VeryLow < Low < Medium < High < VeryHigh
```

```r
# The ordering is NOW imposed
levels(myfac2.temp.data)
```

```
[1] "VeryLow"  "Low"      "Medium"   "High"     "VeryHigh"
```

```r
table(myfac2.temp.data)
```

```
myfac2.temp.data
 VeryLow      Low   Medium     High VeryHigh
```

```
     2        5        2        2        6
```

```
is.factor(myfac2.temp.data)
```

```
[1] TRUE
```

```
is.ordered(myfac2.temp.data)
```

```
[1] TRUE
```

```
# Stripping a factor to the essentials: integer vector
attributes(myfac2.temp.data)
```

```
$levels
[1] "VeryLow"  "Low"       "Medium"   "High"      "VeryHigh"
```

```
$class
[1] "ordered" "factor"
```

```
class(myfac2.temp.data) <- NULL
levels(myfac2.temp.data) <- NULL
myfac2.temp.data
```

```
 [1] 4 2 5 2 1 3 5 5 2 2 3 5 5 5 2 4 1
```

## 5.3   Dates and times in R.

- **Date** class :
  - represents calendar dates
  - built on top of doubles with class attribute 'Date'
  - 0 : Jan 1. 1970 (**Unix Epoch time**)
  - **as.Date()**: method to cast string to a Date
- **POSIXct** and **POSIXlt** : date and time
  - **POSIXct**: stores date/time values as the #seconds since Jan. 1, 1970
  - **POSIXlt**: stored as **bluelist** with elements for seconds, minutes, hours, day, month, year, etc.
- **lubridate**: a very useful package for dates and times:

### 5.3.1   Examples

- **Date**

```
today <- Sys.Date()
today
```

```
[1] "2023-07-25"
```

```r
# Attributes of Date
class(today)
```

```
[1] "Date"
```

```r
attributes(today)
```

```
$class
[1] "Date"
```

```r
unclass(today)
```

```
[1] 19563
```

```r
d0 <- structure(0, class='Date')
d0
```

```
[1] "1970-01-01"
```

```r
class(d0)
```

```
[1] "Date"
```

```r
typeof(d0)
```

```
[1] "double"
```

```r
# Convert a string into a Date
d1  <- as.Date("2022-01-01")
d1
```

```
[1] "2022-01-01"
```

```r
class(d1)
```

```
[1] "Date"
```

```r
typeof(d1)
```

```
[1] "double"
```

- **POSIXct**

```r
# Convert a string into a POSIXct object
now_ct <- as.POSIXct("2018-08-01 22:00", tzone="MST")
now_ct
```

```
[1] "2018-08-01 22:00:00 MDT"
```

```r
attributes(now_ct)
```

```
$class
[1] "POSIXct" "POSIXt"

$tzone
[1] ""
```

```r
typeof(now_ct)
```

```
[1] "double"
```

```r
# Removal of the attributes
attr(now_ct,"tzone") <- NULL
unclass(now_ct)
```

```
[1] 1533182400
```

# Bibliography

Seber G.A.F. & Lee A.J. (2012). Linear Regression Analysis. Wiley Series in Probability and Statistics. Wiley.