

Apunts
Programació 2
Grau en Estadística

Marta Fairén
Edelmira Pasarella

10 d'abril de 2017

Índex

1	Vectors (recordatori)	5
1.1	Creació de vectors al llenguatge R	6
1.2	Accés als components d'un vector	8
1.3	Recorregut	10
1.4	Cerca	13
1.5	Sorting	14
2	Matrius	17
2.1	Creació de matrius al llenguatge R	18
2.2	Accés als components d'una matriu	21
2.3	Recorregut	23
2.4	Cerca	25
3	Lists	29
3.1	Creació de lists	30
3.2	Accés als components d'un list	32
3.3	Afegir i eliminar elements d'un list	35
3.4	Accés a sublists	36
3.5	Recorreguts i cerques	37
4	Data Frames	39
4.1	Creació de data frames al llenguatge R	40
4.1.1	Creació directa	40
4.1.2	Creació de data frames a partir de fitxers	41
4.2	Operadors del tipus <i>list</i> i data frames	41
4.2.1	Accés a una columna	42
4.2.2	Obtenir la longitud	42
4.2.3	Noms	42
4.3	Operadors del tipus matriu i data frames	43
4.3.1	Accés a una columna sencera	44
4.3.2	Accés a una fila sencera	44
4.3.3	Accés a una component d'un data frame	45
4.4	Subdata frames	45
4.4.1	Selecció per rang d'índex	46
4.4.2	Filtres	47
4.4.3	Selecció de files	47
4.5	Extensió de data frames	48
4.5.1	Afegir una fila	48

4.5.2	Afegir una nova columna	50
4.6	Modificació de data frames	51
4.7	Ordenació de data frames	52
4.8	Fusió de data frames	54

Capítol 1

Vectors (recordatori)

Comencem per recordar que un **vector** v és un objecte **estructurat** compostat per un nombre finit d'objectes N anomenats components o elements. Tots els components d'un vector són del mateix tipus T . El tipus T és el **tipus base del vector**. És a dir, els vectors són objectes homogenis. De manera genèrica, el vector v del que estem parlant és de la forma:

t_1	t_2	t_3	t_4	t_5	t_6	\dots	t_N
-------	-------	-------	-------	-------	-------	---------	-------

on $\forall i \in \{1, \dots, N\}$, $t_i \in T$. És a dir, t_i es un valor del tipus base T . Un vector d'enters, de dimensió $N = 8$ podria ser aquest de sota:

1	2	3	4	-4	-3	-2	-1
---	---	---	---	----	----	----	----

Mentre que un vector de booleans, de dimensió $N = 4$ pot ser per exemple:

<i>TRUE</i>	<i>TRUE</i>	<i>FALSE</i>	<i>TRUE</i>
-------------	-------------	--------------	-------------

Pel que fa a requeriment de memòria de les variables de tipus vector, els seus components s'emmagatzemen de manera contigua a un bloc de memòria de mida igual a la dimensió del vector. Per a ser precisos hem d'aclarir que, en realitat, per a una variable vector es reserva espai a memòria equivalent a la dimensió del vector multiplicat per l'espai requerit per emmagatzemar un objecte del tipus base. El tipus base no sempre és un dels tipus base del lleguatge. Això vol dir que no sempre tenim vectors d'enters, booleans, reals i caracters. Podem tenir vectors de tipus estructurats també. En particular, els vector amb tipus base vectors (vectors de vectors) són coneguts com matrius i d'això en parlarem al Capítol 2.

1.1 Creació de vectors al llenguatge R

Al llenguatge R tenim varies funcions que són constructors de vectors.

1. La funció `c()`

Aquest constructor ens permet crear vectors per enumeració dels seus components. El vector resultant tindrà dimensió igual al nombre d'elements que s'han fet servir a la crida de la funció `c()`.

```

1 > v <- c(0.4, 0.1)           # N = 2
2 > v
3 [1] 0.4 0.1
4 >
5 > v <- c(1, 2, 3, 4, 5, 6, 7)   # N = 7
6 > v
7 [1] 1 2 3 4 5 6 7
8 >
9 > v <- c(TRUE, FALSE)         # N = 2
10 > v
11 [1] TRUE FALSE
12 >
13 > v <- c("a", "b", "c", "d", "e") # N = 5
14 > v
15 [1] "a" "b" "c" "d" "e"
16 >

```

S'ha de notar que poden crear vectors buits:

```

1 > v <- c()
2 > v
3 NULL
4 >

```

2. La funció `vector(Tipus_Base, Dimensió)`

Aquest altre constructor ens permet crear vectors d'un tipus base i d'una dimensió donada. Si el `Tipus_Base` no es dona, per defecte, crea un vector de booleans inicialitzat en `FALSE`.

```

1 > v <- vector(length=6)
2 > v
3 [1] FALSE FALSE FALSE FALSE FALSE FALSE
4 >

```

3. L'operador `:`

Aquest operador permet crear vectors d'un rang de números.

```

1 > 1:10
2 [1] 1 2 3 4 5 6 7 8 9 10
3 > v <- 1:10
4 > v[2]

```

```

5 [1] 2
6 > v[2] <- 4
7 > v
8 [1] 1 4 3 4 5 6 7 8 9 10
9 >
10 > v <- 4:9
11 > v
12 [1] 4 5 6 7 8 9
13 >
14 > v <- 12:7
15 > v
16 [1] 12 11 10 9 8 7
17 >

```

4. Altres operadors que ens permeten construir vectors: `seq()` and `rep()`,

```

1 > v <- seq(from=16,to=24, by=2)
2 > v
3 [1] 16 18 20 22 24
4 >
5 > v <- seq(from=16,to=4, by=-3)
6 > v
7 [1] 16 13 10 7 4
8 >
9 > v <- seq(8)
10 > v
11 [1] 1 2 3 4 5 6 7 8
12 >
13 > v <- rep(1,5)
14 > v
15 [1] 1 1 1 1 1
16 >
17 > v <- rep(15,3)
18 > v
19 [1] 15 15 15
20 >

```

I també podem combinar els operadors:

```

1 > v <- c(10,rep(2,5),9)
2 > v
3 [1] 10 2 2 2 2 2 9
4 >
5 > v <- c(seq(5))
6 > v
7 [1] 1 2 3 4 5
8 > v <- c(13:20)
9 > v
10 [1] 13 14 15 16 17 18 19 20
11 >
12 > v <- c(seq(from=15, to=25, by=3))

```

```

13 > v
14 [1] 15 18 21 24
15 >

```

Un cop s'ha creat un vector v , no es pot canviar la seva dimensió. Malgrat això, podem "redimensionar" un vector con segueix:

```

1 > v <- 1:4
2 > v
3 [1] 1 2 3 4
4 > v <- c(v,10)
5 > v
6 [1] 1 2 3 4 10
7 > v <- c(v,10)
8 > v <- c(10,v)
9 > v
10 [1] 10 1 2 3 4 10 10
11 > v <- c(10:20,v)
12 > v
13 [1] 10 11 12 13 14 15 16 17 18 19 20 10 1 2 3 4
    10 10

```

Amb les instruccions de dalt sembla que estem redimensionant el vector v . En realitat, aquestes instruccions són força costoses perquè cada vegada que s'afegeix un o més components, es crea un nou vector a un altre lloc de memòria i el vector original s'ha de relocalitzar totalment. Llavors, s'ha de fer servir amb molta cura aquest tipus d'instruccions.

1.2 Accés als components d'un vector

Mitjançant un **operador d'accés directe**, que té com a paràmetre un enter que correspon a un índex que indica la posició que es desitja visitar, es pot accedir a qualsevol component del vector de manera directa. Matemàticament l'especificació d'aquesta funció és la següent:

$$[\] : \text{vector} \times \text{enter} \rightarrow T$$

Llavors, donat un vector v i una expressió de tipus enter exp , fem servir la crida

$$v[exp]$$

amb

$$1 \leq \text{avaluacio}(exp) \leq N \quad (1)$$

per referir-nos al component del vector que es troba a la posició que resulti d'avaluar exp .

Aleshores, no cal visitar els components anteriors a un component donat per arribar-hi.

```

1 > v <- c(0.4, 0.1)
2 > v[1]
3 [1] 0.4

```



```

4 >
5 > v <- c(1, 2, 3, 4, 5, 6, 7)
6 > v[5]
7 [1] 5
8 >
9 > v <- c(TRUE, FALSE)
10 > v[2]
11 [1] FALSE
12 > v <- c("a", "b", "c", "d", "e")
13 > i <- 2
14 > v[i]
15 [1] "b"
16 > v[i+1]
17 [1] "c"
18 > v[i*2]
19 [1] "d"
20 > j <- 3
21 > v[i+j]
22 [1] "e"
23 > v[2*i+1]
24 [1] "e"
25 >

```

Si assignem un valor numèric a qualsevol component d'un vector creat amb l'operador `vector()`, el tipus base canvia.

```

1 > v<-vector(length=6)
2 > v
3 [1] FALSE FALSE FALSE FALSE FALSE FALSE
4 >
5 > v[4] <- 55
6 > v
7 [1] 0 0 0 55 0 0
8 >

```

És molt important vigilar que es satisfaci la restricció (1). En cas que no la respectem tindrem un problema en intentar accedir a una posició no definida dins del vector. En aquests casos, l'R ens respon amb `NA` (Not Available):

```

1 > v <- c("a", "b", "c", "d", "e")
2 > v[6]
3 [1] NA
4 >

```

En particular, quan el vector és buit, tindrem:

```

1 > v <- c()
2 > v
3 NULL
4 > v[1]
5 NULL
6 >

```

S'ha de tenir en compte que, segons la definició de vectors, $v[exp] = t_{avaluacio(exp)} \in T$ i que per tant, aquest component es comporta com una variable del tipus base i podem fer servir qualsevol dels seus operadors.

```

1 > w <- rep(8,7)
2 > w
3 [1] 8 8 8 8 8 8 8
4 > v <- 5:10
5 > v
6 [1] 5 6 7 8 9 10
7 > cat(v[5]+w[6], "\n")
8 17
9 > y <- 5*v[3]
10 > y
11 [1] 35
12 >

```

Ara que hem vist la necessitat que els accessos als components d'un vector siguin posicions permeses, surt la necessitat de conèixer la dimensió d'un vector. Per fer-ho disposem a l'R d'una funció especificada com segueix:

$$\text{length} : \text{vector} \rightarrow \text{enter}$$

Així doncs, si la dimensió d'un vector v és N , la crida

$$\text{length}(v)$$

ens tornarà N .

```

1 > v<-c()
2 > length(v)
3 [1] 0
4 >
5 > v <- c("a", "b", "c", "d", "e")
6 > length(v)
7 [1] 5
8 >
9 > v <- c(TRUE, FALSE)
10 > length(v)
11 [1] 2
12 >

```

1.3 Recorregut

Hi ha molts problemes on s'han de processar tots i cadascun dels components d'una seqüència que, aleshores, pot estar representada per un vector. Aquesta classe de problemes se'ls anomena "Recorreguts". Per exemple, per calcular la mitjana d'una seqüència de notes d'estudiants, per avaluar un polinomi donat els seus coeficients, hem de fer un recorregut. Hi ha molts problemes que requereixen que tots els elements de l'entrada es processin. Llavors, per fer recorreguts d'un vector, tindrem un esquema genèric com aquest:

```

1 A
2 for (i in 1:length(v)){
3   processar(v[i])
4 }
5 B

```

On A correspon a un bloc d'instruccions de preprocessament, per exemple, la inicialització d'algunes variables, i B correspon a un bloc d'instruccions de postprocessament, per exemple, imprimir el resultat.

Un primer exemple molt simple és l'escriptura dels components d'un vector:

```

1 escriure_vector <- function(v){
2   for (i in 1:length(v)){
3     cat(v[i],"\n")
4   }
5 }

```

I es pot utilitzar d'aquesta forma:

```

1 > v<- 2:15
2 > escriure_vector(v)
3 2
4 3
5 4
6 5
7 6
8 7
9 8
10 9
11 10
12 11
13 12
14 13
15 14
16 15
17 >

```

Suposem que necessitem una funció per llegir un vector. Aquesta funció pot ser con segueix

```

1 llegir_vector <- function(v){
2   for (i in 1:length(v)){
3     v[i] <- scan(n=1,quiet=TRUE)
4   }
5   return (v)
6 }

```

I la podem fer servir de la següent manera:

```

1 > v <- vector(length = 8)
2 > v <- llegir_vector(v)
3 1: 2
4 1: 4
5 1: 6

```

```

6 1: 8
7 1: 10
8 1: 12
9 1: 14
10 1: 16
11 > v
12 [1] 2 4 6 8 10 12 14 16
13 >

```

Amdós exemples previs, no requereixen cap preprocessament o postprocesament. Considerem ara el problema de trobar el valor màxim dins d'un vector no buit (aquesta és la preconditionió del problema):

```

1 maxim <- function(v){
2   max<- v[1]
3   for(i in 2:length(v)){
4     if (v[i] > max) max <- v[i]
5   }
6   return (max)
7 }

1 > v<-c(12,5,6,7,10,100,4,5, 56,2,20)
2 > maxim(v)
3 [1] 100
4 > v<-c(12,5,6,7,10,100,4,5, 56,2,200)
5 > maxim(v)
6 [1] 200
7 > v<-c(1200,5,6,7,10,100,4,5, 56,2,20)
8 > maxim(v)
9 [1] 1200
10 >

```

Un altre exemple és la funció que calcula el producte escalar de dos vectors de la mateixa dimensió (preconditionió):

```

1 producte_escalar <- function(v,w){
2   p <- 0
3   for (i in 1:length(v)){
4     p <- p + v[i]*w[i]
5   }
6   return (p)
7 }

1 > v <- rep(2,3)
2 > v
3 [1] 2 2 2
4 > w <- 1:3
5 > w
6 [1] 1 2 3
7 > producte_escalar(v,w)
8 [1] 12
9 >

```

1.4 Cerca

Hi ha una classe de problemes on no cal visitar o processar tots els elements, és a dir, fer un recorregut, perquè ens demanen confirmar un predicat, és a dir, una propietat, $P(s)$, sobre la seqüència s representada amb un vector v , com per exemple:

- $P(s)$ = "la seqüència s té múltiples de 3".

Aquesta propietat vol dir el següent: $\exists i : 1 \leq i \leq n : v[i] \bmod 3 = 0$ i això s'ha de comprovar. Anirem processant cadascun dels components del vector i si trobem un que sigui múltiple de 3, hem d'aturar la comprovació i donar una resposta afirmativa. En cas que no hi hagi cap múltiple de 3 al vector, acabarem fent un recorregut i la resposta seria negativa. Un possible script a l'R és el següent:

```

1 multiploe3 <- function(v){
2   i <- 1
3   trobat <- FALSE
4   while (i <= length(v) && !trobat){
5     trobat <- v[i] %% 3 == 0
6     i <- i + 1
7   }
8   return (trobat)
9 }
```

- $P(s)$ = "la seqüència s té un cim".

Per a aquest problema és necessari conèixer la definició de cim. Atès que la seqüència s es representa com un vector s , un cim és un component $v[i]$, $2 \leq i \leq n - 1$, de v tal que $v[i - 1] < v[i] < v[i + 1]$. Si, mentre es fa el processament del vector es troba un component que satisfaci aquesta propietat s'ha d'aturar el recorregut perquè ja podem donar una resposta afirmativa. En cas contrari, si haguéssim de fer el processament de tots els components del vector, és a dir, el recorregut del vector, obtindríem una resposta negativa. Obviament, per comprobar si un vector té un cim, es requereix que tingui com a mínim tres elements. Llavors, aquesta és la precondició del problema: la dimensió de v ha de ser més gran o igual a 3. Un possible script a l'R és el següent:

```

1 #Precondicio: length(v) >= 3
2 cim <- function(v){
3   ant <- 1
4   act <- 2
5   suc <- 3
6   trobat <- FALSE
7   while (suc <= length(v) && !trobat){
8     trobat <- v[ant] < v[act] && v[act] > v[suc]
9     ant <- act
10    act <- suc
11    suc <- suc + 1
12  }
13  return (trobat)
14 }
```

- $P(\mathbf{s}) =$ "la seqüència \mathbf{s} és creixent"

Si analitzem aquesta propietat, ens adonarem que en realitat vol dir el següent: $\forall i : 1 \leq i \leq n - 1 : v[i] \leq v[i + 1]$ i això és el que realment hem de comprovar. En aquest cas, la cerca consisteix en trobar un contraexemple d'aquesta propietat. És a dir, si trobem un parell de components consecutius del vector tal que $1 \leq i \leq n - 1 : v[i] > v[i + 1]$ hem d'aturar la comprovació de la propietat i donar una resposta negativa. En aquest cas, si féssim el recorregut, la resposta seria afirmativa. S'ha de notar que, segons la definició, tota seqüència buida o que tingui un sol element és creixent. Un possible script en R és el següent:

```

1 creixent <- function(v){
2   i <- 1
3   continuar <- TRUE
4   while (i < length(v) && continuar){
5     continuar <- v[i] < v[i+1]
6     i <- i + 1
7   }
8   return (continuar)
9 }
```

- $P(\mathbf{s}) =$ "la seqüència \mathbf{s} no té cap element parell".

Si analitzem aquesta propietat, ens adonarem que en realitat vol dir el següent: $\forall i : 1 \leq i \leq n : v[i] \bmod 2 \neq 0$ no és parell i això és el que realment hem de comprovar. Un altre cop, el problema consisteix en trobar un contraexemple de la propietat. És a dir, si trobem un component parell del vector que representa la seqüència, hem d'aturar la comprovació i donar una resposta negativa. Si féssim el recorregut, la resposta seria afirmativa. Un possible script a l'R és el següent:

```

1 noparells <- function(v){
2   i <- 1
3   continuar <- TRUE
4   while (i <= length(v) && continuar){
5     continuar <- v[i] %% 2 != 0
6     i <- i + 1
7   }
8   return (continuar)
9 }
```

1.5 Sorting

Una família molt important d'algorismes és la família d'algorismes d'ordenament d'un vector. El problema pot ser vist de la següent manera: donat un vector, s'han de reorganitzar els seus components de forma tal que el vector arribi a representar una seqüència creixent (també pot ser decreixent). Això es pot formalitzar com segueix:

$$\text{sort}(v, v') = \text{permutacio}(v, v') \text{ i } \forall i : 1 \leq i \leq n - 1 : v'[i] \leq v'[i + 1]$$

On el predicat $permutacio(v, v')$ diu que v' és una permutació de v . Els algorismes bàsics d'ordenament són el `SelectSort` (ordenament per selecció), l'`InsertSort` (ordenament per inserció) i el `BubbleSort` (ordenament per intercanvi o mètode de la bombolla). La dificultat que tenen aquests algorismes és que el vector s'ha d'ordenar sense fer servir cap emmagatzament addicional. És a dir, només s'ha de fer servir l'espai a memòria del vector v . Podeu trobar més detalls d'aquests algorismes a [1]. A sota presentem una implementació a l'R del `SelectSort`.

```

1  buscarmin <- function(v, inicio, final){
2    pos_min <- inicio
3    i <- inicio+1
4    while (i <= final){
5      if (v[i] < v[pos_min]) pos_min <- i
6      i <- i+1
7    }
8    return (pos_min)
9  }
10
11 SelectSort <- function(v){
12   n <- length(v)
13   i <- 1
14   #Durant la resta de l'execucio es mante que
15   #el segment v[1..i-1] esta ordenat
16   while (i <= n-1){
17     #Busca el component minim a v[i..n]
18     pos_min <- buscarmin(v,i,n)
19     #Intercanvia v[i] amb v[pos_min]
20     aux <- v[i]
21     v[i] <- v[pos_min]
22     v[pos_min] <- aux
23     i <- i+1
24   }
25   return (v)
26 }

```

La funció `sort` a l'R. El llenguatge R disposa d'una funció que ordena un vector de manera creixent (opció per defecte) o decreixent. La discussió en relació a l'algorisme d'ordenament que implementa aquesta funció està fora de l'abast d'aquest document.

```

1 > v <- c(9,3,5,1,7,9,5) sort(v) [1] 1 3 5 5 7 9 9
2 > sort(v,decreasing=TRUE) [1] 9 9 7 5 5 3 1

```


Capítol 2

Matrius

Una **matriu** m és un objecte **estructurat**, bidimensional, que pot ser vist com una taula. Des del punt de vista d'estructures de dades, una matriu és un vector de vectors. És a dir, un vector on cada component és un vector. Aleshores, al igual que els vectors, tots els components d'una matriu són del mateix tipus base T . Llavors, una matriu és un objecte homogeni, compost per un nombre finit de files F (entrades horitzontals) i un nombre finit de columnes C (entrades verticals) anomenats components o elements. De manera genèrica, la matriu $m_{F \times C}$ de la que estem parlant és de la forma:

$$\begin{array}{c} \overbrace{\hspace{10em}}^{C \text{ columnes}} \\ \left[\begin{array}{ccccc} m_{1,1} & m_{1,2} & \dots & m_{1,C-1} & m_{1,C} \\ m_{2,1} & m_{2,2} & \dots & m_{2,C-1} & m_{2,C} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ m_{F-1,1} & m_{F-1,2} & \dots & m_{F-1,C-1} & m_{F-1,C} \\ m_{F,1} & m_{F,2} & \dots & m_{F,C-1} & m_{F,C} \end{array} \right] \\ \underbrace{\hspace{10em}}_{F \text{ files}} \end{array}$$

Com podem veure, la matriu m té dimensió $F \times C$ i per tant té $F \times C$ components on $\forall i \in \{1, \dots, F\} \forall j \in \{1, \dots, C\}, m_{i,j} \in T$. És a dir, $m_{i,j}$ és un valor del tipus base T . Una matriu d'enters, de dimensió 4×4 podria ser aquesta:

$$m_{4 \times 4} = \begin{bmatrix} 1 & 2 & 3 & 7 \\ 4 & 4 & 3 & 2 \\ 2 & 1 & 6 & 1 \\ 3 & 1 & 5 & 1 \end{bmatrix}$$

Mentre que una matriu de booleans, de dimensió 3×2 pot ser per exemple:

$$m_{3 \times 2} = \begin{bmatrix} TRUE & TRUE \\ FALSE & TRUE \\ TRUE & FALSE \end{bmatrix}$$

2.1 Creació de matrius al llenguatge R

Al llenguatge R tenim la següent funció per construir matrius.

Funció `matrix()`

Aquest constructor ens permet crear matrius per enumeració dels seus components, i.e. mitjançant un vector.

1. El nombre de files, $matrix(c(e_1, \dots, e_n), nrow = F)$

```
1 > m <- matrix(c(1,2,3,4,5,6,7,8,9,10, 11,12), nrow
  =4)
2 > m
3      [,1] [,2] [,3]
4 [1,]    1    5    9
5 [2,]    2    6   10
6 [3,]    3    7   11
7 [4,]    4    8   12
8 >
```

2. El nombre de columnas, $matrix(c(e_1, \dots, e_n), ncol = C)$

```
1 > m <- matrix(c(1,2,3,4,5,6,7,8,9,10, 11,12), ncol
  =4)
2 > m
3      [,1] [,2] [,3] [,4]
4 [1,]    1    4    7   10
5 [2,]    2    5    8   11
6 [3,]    3    6    9   12
7 >
```

3. Ambdós, $matrix(c(e_1, \dots, e_n), nrow = F, ncol = C)$

```
1 > m <- matrix(c(1,2,3,4,5,6,7,8,9,10, 11,12), nrow
  =4, ncol=3)
2 > m
3      [,1] [,2] [,3]
4 [1,]    1    5    9
5 [2,]    2    6   10
6 [3,]    3    7   11
7 [4,]    4    8   12
8 >
```

4. Cap modalitat $matrix(c(e_1, \dots, e_n))$. En aquest cas els elements es col·loquen en una matriu d'una única columna.

```

1 > m <- matrix(c(1,2,3,4,5,6,7,8,9,10, 11,12))
2 > m
3      [,1]
4 [1,]    1
5 [2,]    2
6 [3,]    3
7 [4,]    4
8 [5,]    5
9 [6,]    6
10 [7,]    7
11 [8,]    8
12 [9,]    9
13 [10,]   10
14 [11,]   11
15 [12,]   12
16 >

```

IMPORTANT Els components d'una matriu creada amb la funció `matrix()` són emmagatzemats per columnes.

La matriu resultant tindrà dimensió igual al nombre d'elements que s'han fet servir a la crida de la funció `matrix()`, és a dir n . Aleshores, n ha de ser múltiple de les files i de les columnes.

```

1 > m <- matrix(c(1,2,3,4,5,6,7,8,9,10,11,12), nrow=5)
Mensajes de aviso perdidos
In matrix(c(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12), nrow = 5) :
la longitud de los datos [12] no es un submúltiplo o múltiplo del número
de filas [5] en la matriz
1 > m <- matrix(c(1,2,3,4,5,6,7,8,9,10, 11,12), ncol=7)
Mensajes de aviso perdidos
In matrix(c(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12), ncol = 7) :
la longitud de los datos [12] no es un submúltiplo o múltiplo del número
de columnas [7] en la matriz
1 > m <- matrix(c(1,2,3,4,5,6,7,8,9,10, 11,12), nrow=2)
2 > m
3      [,1] [,2] [,3] [,4] [,5] [,6]
4 [1,]    1    3    5    7    9   11
5 [2,]    2    4    6    8   10   12
6 >

```

De tota manera, si volem emmagatzemar els components per files podem fer servir la modalitat `matrix(c(...),nrow=F, byrow=TRUE)`

Com al cas dels vectors, es poden crear matrius "buides". En realitat són matrius on els seus components tenen el valor NA.

```

1 > m <- matrix(nrow=4,ncol=3)
2 > m
3      [,1] [,2] [,3]
4 [1,]   NA   NA   NA

```

```

5 [2,] NA NA NA
6 [3,] NA NA NA
7 [4,] NA NA NA
8 > m <- matrix(nrow=4)
9 > m
10      [,1]
11 [1,] NA
12 [2,] NA
13 [3,] NA
14 [4,] NA
15 > m <- matrix(ncol=3)
16 > m
17      [,1] [,2] [,3]
18 [1,] NA NA NA
19 >
20 > m <- matrix()
21 > m
22      [,1]
23 [1,] NA
24 >

```

A diferència dels vectors, després d'assignar un valor d'un tipus T a algú dels components de la matriu, la resta dels components mantenen el valor NA.

```

1 > m <- matrix(nrow=4,ncol=3)
2 > m
3      [,1] [,2] [,3]
4 [1,] NA NA NA
5 [2,] NA NA NA
6 [3,] NA NA NA
7 [4,] NA NA NA
8 > m[3,2] <- 100
9 > m
10      [,1] [,2] [,3]
11 [1,] NA NA NA
12 [2,] NA NA NA
13 [3,] NA 100 NA
14 [4,] NA NA NA
15 > m <- matrix(nrow=4,ncol=3)
16 > m[3,2] <- TRUE
17 > m
18      [,1] [,2] [,3]
19 [1,] NA NA NA
20 [2,] NA NA NA
21 [3,] NA TRUE NA
22 [4,] NA NA NA
23 >

```

2.2 Accés als components d'una matriu

S'accedeix als components d'una matriu mitjançant un **operador d'accés directe**, que té com a paràmetres dos enters que corresponen a un índex per a la fila i un índex per a la columna que es desitja accedir, respectivament. És a dir, indicant la posició que es desitja visitar es pot accedir a qualsevol component del vector de manera directa. Matemàticament l'especificació d'aquesta funció és la següent:

$$[\] : \text{matriu} \times \text{enter} \times \text{enter} \rightarrow T$$

Llavors, donada una matriu m i un parell d'expressions de tipus enter, exp_f i exp_c , fem servir la crida

$$m[exp_f, exp_c]$$

amb

$$1 \leq \text{avaluacio}(exp_f) \leq F \quad (2)$$

$$1 \leq \text{avaluacio}(exp_c) \leq C \quad (3)$$

per referir-nos al component de la matriu que es troba a la posició $\langle exp_f, exp_c \rangle$.

Aleshores, com al cas dels vectors, tenim accés directe per visitar els components d'una matriu.

```

1 > m
2      [,1] [,2] [,3]
3 [1,]    1    5    9
4 [2,]    2    6   10
5 [3,]    3    7   11
6 [4,]    4    8   12
7 > m[3,2]
8 [1] 7
9 > m[1,3]
10 [1] 9
11 > i <- 2
12 > m[2*i,i]
13 [1] 8
14 > j <- 2
15 > m[i+j,3]
16 [1] 12
17 >

```

Una de les facilitats per manipular matrius que ens ofereix l'R és que podem accedir no només als components de les matrius si no també a submatrius. En particular, podem agafar una fila o una columna sencera.

```

1 > m
2      [,1] [,2]
3 [1,]    1    4
4 [2,]    2    5
5 [3,]    3    6
6 > m[,2]      #accedeix a la columna 2
7 [1] 4 5 6
8 > m[2,]      #accedeix a la fila 2
9 [1] 2 5

```

```

10 > m[1:2,]
11      [,1] [,2]
12 [1,]    1    4
13 [2,]    2    5
14 >

```

És molt important vigilar que es satisfacin les restriccions (2) i (3). En cas que no es respectin aquestes restriccions, rebrem un missatge d'error en intentar accedir a una posició il·legal dins de la matriu considerada.

```

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

```

Error en m[3, 7] : subíndice fuera de los límites

Semblant als vectors, els components d'una matriu pertanyen al tipus base T , llavors $m[exp_f, exp_c] \in T$ i que per tant, aquest component es comporta com una variable del tipus T i podem fer servir qualsevol dels seus operadors.

```

1 > b
2      [,1] [,2]
3 [1,] TRUE  TRUE
4 [2,] FALSE TRUE
5 [3,] FALSE FALSE
6 [4,] TRUE  FALSE
7 > b[1,2] && b[3,1]
8 [1] FALSE
9 > b[1,2] && b[4,1]
10 [1] TRUE
11 > m
12      [,1] [,2]
13 [1,]    1    4
14 [2,]    2    5
15 [3,]    3    6
16 > cat(m[2,2]+m[3,1], "\n")
17 8
18 >

```

A més dels operadors mencionats a dalt, podem conèixer la dimensió d'una matriu i el número de components que té mitjançant els següents operadors:

dim : matriu \rightarrow enter \times enter

length : matriu \rightarrow enter

nrow : matriu \times enter \rightarrow enter

ncol : matriu \times enter \rightarrow enter

```

1 > m <- matrix(c(1,2,3,4,5,6),ncol=2)
2 > m
3      [,1] [,2]
4 [1,]    1    4
5 [2,]    2    5
6 [3,]    3    6
7 > dim(m)
8 [1] 3 2
9 > c <- dim(m)
10 > c
11 [1] 3 2
12 > length(m)
13 [1] 6
14 > nrow(m)
15 [1] 3
16 > ncol(m)
17 [1] 2
18 >

```

2.3 Recorregut

Hi ha molts problemes on s'han de processar tots i cadascun dels components d'una matriu. Com ja hem dit al capítol 1, aquesta classe de problemes se'ls anomena "Recorreguts". Per recórrer tots els elements d'una matriu hem d'accedir a tots ells. Això ho podem fer per files o per columnes. El següent esquema és un recorregut per files d'una matriu:

```

1 A
2 for (i in 1:nrow(m)){
3   for (j in 1:ncol(m)){
4     processar(m[i,j])
5   }
6 }
7 B

```

On A correspon a un bloc d'instruccions de preprocessament, per exemple, la inicialització d'algunes variables, i B correspon a un bloc d'instruccions de postprocessament, per exemple, imprimir el resultat.

També es pot fer el recorregut d'una matriu per columnes amb l'esquema de sota:

```

1 A
2 for (j in 1:ncol(m)){
3   for (i in 1:nrow(m)){
4     processar(m[i,j])
5   }
6 }
7 B

```

Com a exemple podem veure l'escriptura dels components d'una matriu m :

```

1 escriure_matriu <- function(m){
2   for (i in 1:nrow(m)){
3     for (j in 1:ncol(m)){
4       cat(m[i,j], "\n")
5     }
6   }
7 }

```

I aquesta funció es pot utilitzar d'aquesta forma:

```

1 > m
2      [,1] [,2]
3 [1,]    1    2
4 [2,]    3    4
5 [3,]    5    6
6 > escriure_matriu(m)
7 1
8 2
9 3
10 4
11 5
12 6
13 >

```

També podem donar format de matriu a aquesta sortida fent la funció com segueix:

```

1 escriure_matriu <- function(m){
2   for (i in 1:nrow(m)){
3     for (j in 1:ncol(m)){
4       cat(m[i,j], " ")
5     }
6       cat(m[i,j], "\n")
7   }
8 }

```

I això donarà la sortida següent:

```

1 > m
2      [,1] [,2]
3 [1,]    1    2
4 [2,]    3    4
5 [3,]    5    6
6 > escriure_matriu(m)
7 1 2
8 3 4
9 5 6
10 >

```

Suposem que necessitem una funció per llegir (per files) una matriu. Aquesta funció pot ser con segueix

```

1 llegir_matriu <- function(F,C){
2   m <- matrix(nrow=F,ncol=C)

```



```

3   for (i in 1:F){
4     for (j in 1:C){
5       m[i,j] <- scan(n=1,quiet=TRUE)
6     }
7   }
8   return (m)
9 }

```

I la podem fer servir de la següent manera:

```

1 > m <- llegir_matriu(3,2)
2 1: 1
3 1: 2
4 1: 3
5 1: 4
6 1: 5
7 1: 6
8 > m
9      [,1] [,2]
10 [1,]    1    2
11 [2,]    3    4
12 [3,]    5    6
13 >

```

2.4 Cerca

Com al cas dels vectors, tenim una classe de problemes on no cal visitar o processar tots els elements, és a dir, fer un recorregut, perquè ens demanen confirmar un predicat, és a dir, una propietat, $P(m)$, sobre la matriu m . Per resoldre els problemes de cerca, s'han de fer servir les mateixes tècniques que hem explicat en la Secció 1.4. Per exemple:

$P(s) =$ "els components de la matriu $m_{F \times C}$ formen una seqüència estrictament creixent quan es recorren per fila".

Si analitzem aquesta propietat, ens adonarem que en realitat vol dir el següent: $\forall i \forall j : 1 \leq i \leq F : 1 < j \leq C : m[i, (j-1)] \leq m[i, j]$ i $\forall i : 1 < i \leq F : m[(i-1), C] \leq m[i, 1]$. Llavors, això és el que realment hem de comprovar. En aquest cas, la cerca consisteix en trobar un contraexemple d'aquesta propietat. És a dir, si trobem un parell de components consecutius del recorregut per files de la matriu tal que el primer sigui més gran que el segon hem d'aturar la comprovació de la propietat i donar una resposta negativa. S'ha de notar que, segons la definició, tota seqüència buida o que tingui un sol element és creixent. Un possible script en R és el següent:

```

1 creixent <- function(m){
2   ant <- m[1,1]
3   for (i in 2:nrow(m)){
4     for (j in 1:ncol(m)){
5       if (m[i,j] <= ant) return (FALSE)
6       else ant <- m[i,j]

```

```
7     }
8   }
9   return (TRUE)
10 }
11
12 > m <- matrix(c(1,2,3,4,5,6,7,8,9),nrow=3,byrow=TRUE)
13 > m
14     [,1] [,2] [,3]
15 [1,]    1    2    3
16 [2,]    4    5    6
17 [3,]    7    8    9
18 >
19 > creixent(m)
20 [1] TRUE
21 > m <- matrix(c(1,2,3,3,4,5,5,5,6),nrow=3,byrow=TRUE)
22 > m
23     [,1] [,2] [,3]
24 [1,]    1    2    3
25 [2,]    3    4    5
26 [3,]    5    5    6
27 > creixent(m)
28 [1] FALSE
29 >
30 > m <- matrix(c(1,2,3,3,4,5,5,5,6),nrow=3)
31 > m
32     [,1] [,2] [,3]
33 [1,]    1    3    5
34 [2,]    2    4    5
35 [3,]    3    5    6
36 > creixent(m)
37 [1] FALSE
38 >
39 > m <- matrix(c(3,5,6,7,9,11),nrow=3,byrow=TRUE)
40 > m
41     [,1] [,2] [,3]
42 [1,]    9    2    3
43 [2,]    3    4    5
44 [3,]    5    5    6
45 > creixent(m)
46 [1] TRUE
47 >
48 > m <- matrix(c(1,2,3,3,4,5,5,5,-1),nrow=3,byrow=TRUE)
49 > m
50     [,1] [,2] [,3]
51 [1,]    1    2    3
52 [2,]    3    4    5
53 [3,]    5    5   -1
54 > creixent(m)
55 [1] FALSE
56 >
```

Deixem com a exercici al lector fer la implementació corresponent a verificar si una matriu m és estrictament creixent quan es recorre (a) per columnes, (b) en ziga-zaga horitzontal (començant des de la posició $m[1, 1]$ d'esquerra a dreta) i (c) en ziga-zaga vertical (començant des de la posició $m[1, 1]$ d'adalt a baix).

Capítol 3

Lists

Un **list** l és un objecte **estructurat** compost per objectes anomenats components o elements. A diferència d'un vector, els components d'un list poden ser de tipus diferents. Aleshores, els lists són objectes heterogenis. De manera genèrica, un list l és de la forma:

t_1	t_2	t_3	t_4	t_5	t_6	\dots	t_N
-------	-------	-------	-------	-------	-------	---------	-------

on $\forall i \in \{1, \dots, N\}$, $t_i \in T_i$. És a dir, t_i és un valor del tipus T_i i $\forall i \neq j$, $i, j \in \{1, \dots, N\}$ no necessàriament $T_i = T_j$.

Per exemple, considerem el següent list:

<i>Enric</i>	<i>Canals</i>	21	<i>H</i>	1.77	<i>FALSE</i>	66543211	93 - 4444446
--------------	---------------	----	----------	------	--------------	----------	--------------

En aquest exemple el list té components de diferents tipus: cadena de caràcters, enter, real i booleà.

Els components poden rebre un nom que facilita l'accés als components i permet un nivell d'abstracció més alt.

nom_1	nom_2	nom_3	nom_4	nom_5	nom_6	\dots	nom_N
t_1	t_2	t_3	t_4	t_5	t_6	\dots	t_N

En aquest list, $\forall i \in \{1, \dots, N\}$, $t_i \in T_i$ i nom_i és l'identificador del i -èssim component. Com veurem a l'apartat 3.1, els noms d'un list poden ser dinàmicament assignats. A sota podem veure el list de l'exemple anterior amb noms per als components:

<i>Nom</i>	<i>Cognom</i>	<i>Edat</i>	<i>Sexe</i>	<i>Alçada</i>	<i>Treballa</i>	<i>Mobil</i>	<i>Telefon</i>
<i>Enric</i>	<i>Canals</i>	21	<i>H</i>	1.77	<i>FALSE</i>	66543211	93 - 4444446

En aquest exemple el list té components anomenats “Nom” i “Cognom” de tipus cadena de caràcters, “Edat” de tipus enter, “Sexe” de tipus caràcter ('H'/'D'), “Alçada” de tipus real, “Treballa” de tipus booleà (TRUE/FALSE) i, “Mobil” i “Telefon” que podrien ser de tipus cadena de caràcters també.

3.1 Creació de lists

Per a construir un list tenim la següent funció constructora:

```
list()
```

Aquest constructor ens permet crear lists buits, lists per enumeració dels seus components i lists amb *tags* o noms per als seus components, com segueix:

1. list buit

```
1 > l <- list()
2 > l
3 list()
4 >
```

2. list per enumeració

```
1 > l <- list("Enric", "Canals", 21, "H", 1.77, FALSE
2           , 66543211, "93 - 44444446")
3 > l
4 [[1]]
5 [1] "Enric"
6
7 [[2]]
8 [1] "Canals"
9
10 [[3]]
11 [1] 21
12
13 [[4]]
14 [1] "H"
15
16 [[5]]
17 [1] 1.77
18
19 [[6]]
20 [1] FALSE
21
22 [[7]]
23 [1] 66543211
24
25 [[8]]
26 [1] "93 - 44444446"
27 >
```

3. list amb noms per als components

```

1 > l <- list(Nom="Enric", Cognom="Canals",Edat=21,
             Sexe="H",Alcada=1.77,Treballa=FALSE,Mobil
             =66543211,Telefon="93 - 4444446")
2 > l
3 $Nom
4 [1] "Enric"
5
6 $Cognom
7 [1] "Canals"
8
9 $Edat
10 [1] 21
11
12 $Sexe
13 [1] "H"
14
15 $Alcada
16 [1] 1.77
17
18 $Treballa
19 [1] FALSE
20
21 $Mobil
22 [1] 66543211
23
24 $Telefon
25 [1] "93 - 4444446"

```

Adicionalment, podem consultar els noms dels components d'un list mitjançant la funció `names(l)` que donat un list `l` ens retorna un vector amb els noms dels seus components:

```

1 > names(l)
2 [1] "Nom"      "Cognom"    "Edat"     "Sexe"
3 "Alcada"   "Treballa"  "Mobil"    "Telefon"

```

De vegades és necessari donar-li noms als components d'un list creat prèviament per enumeració. Això es pot fer assignant a `names(l)` el vector amb els noms que es vol:

```

1 > l<-list("Maria Perez","FIB",2011)
2 > l
3 [[1]]
4 [1] "Maria Perez"
5
6 [[2]]
7 [1] "FIB"
8
9 [[3]]
10 [1] 2011

```

```

11
12 #Fins aqui els components del list l no tenen
    noms. Ara li afegim els noms:
13
14 > names(l) <- c("Nom", "Centre", "Ingres")
15 > l
16 $Nom
17 [1] "Maria Perez"
18
19 $Centre
20 [1] "FIB"
21
22 $Ingres
23 [1] 2011

```

3.2 Accés als components d'un list

Existeixen diferents maneres d'accedir directament a cada component d'un list. En efecte, mitjançant **operadors d'accés directe**, que tenen com a paràmetre o bé el nom d'un component o bé un índex que indica la posició del component que es desitja visitar, es pot accedir a qualsevol component del list de manera directa. Les especificacions d'aquests operadors són les següents:

1. $\$: list \times id \rightarrow T_{id}$

Essent l un list, id el nom d'un component de l i T_{id} el tipus del component amb nom id . La crida a aquest operador té la següent sintaxi

$$l\$id$$

```

1 > l<-list(Nom="Maria Perez",Centre="FIB",Ingres
    =2011)
2 > l$Nom
3 [1] "Maria Perez"
4 > l$Centre
5 [1] "FIB"
6 > l$Ingres
7 [1] 2011
8 >

```

2. $[[]]: list \times enter \rightarrow T$

En aquest operador d'accés, els lists són tractats de manera semblant als vectors i fem servir la crida

$$l[[exp]]$$

per referir-nos al component del list l que es troba a la posició que resulti d'avaluar l'expressió entera exp . Cal notar, però, una diferència important

amb els vectors, i és que amb el list, l'operador és un doble gafet (`[[]]`) i cal parar atenció a aquesta diferència.

De la mateixa manera que passava amb els vectors, els accessos als components d'un list han de ser necessàriament posicions permeses en el list, i per tant, cal tenir la possibilitat de conèixer la dimensió d'un list. Per fer-ho disposem d'una funció especificada com segueix:

$$\text{length} : \text{list} \rightarrow \text{enter}$$

Així doncs, si un list l té N components, la crida

$$\text{length}(l)$$

ens tornarà N .

```

1 > l<-list(Nom="Maria Perez",Centre="FIB",Ingres
      =2011)
2 > l
3 $Nom
4 [1] "Maria Perez"
5
6 $Centre
7 [1] "FIB"
8
9 $Ingres
10 [1] 2011
11
12 > length(l)
13 [1] 3
14 >
15 > l1 <-list()
16 > length(l1)
17 [1] 0

```

Ara podem demanar que

$$1 \leq \text{avaluacio}(\text{exp}) \leq \text{length}(l) \quad (2)$$

```

1 > l<-list(Nom="Maria Perez",Centre="FIB",Ingres
      =2011)
2 > l[[1]]
3 [1] "Maria Perez"
4 > l[[2]]
5 [1] "FIB"
6 > l[[3]]
7 [1] 2011
8 > l[[1+1]]
9 [1] "FIB"
10 > l[[2*1]]
11 [1] "FIB"

```

Com en el cas dels vectors, és molt important vigilar que es satisfaci la restricció (2). En cas que no la respectem tindrem un problema en intentar accedir a una posició il·legal dins del list. En aquests casos, l'R ens respon con segueix

```
1 > l<-list(Nom="Maria Perez",Centre="FIB",Ingres
      =2011)
2 > l[[4]]
3 Error en l[[4]] : subindice fuera de los limites
4 >
```

S'ha de tenir en compte que existeix una bijecció del conjunt de noms d'un list l i la posició en que aquest nom apareix al list. Per exemple, al list l obtingut fent la instrucció

```
1 l<-list(Nom="Maria Perez",Centre="FIB",Ingres
      =2011):
```

tindríem la bijecció següent:

	<i>posicio</i>
<i>Nom</i>	→ 1
<i>Centre</i>	→ 2
<i>Ingres</i>	→ 3

3. Entre les dues opcions anteriors per accedir als components d'un list l , també existeix una altra manera d'indexar-lo amb la crida següent:

```
l[[quote(id)]]
```

on *quote(id)* és el nom corresponent al *id* dins de doble cometes:

```
1 > l<-list(Nom="Maria Perez",Centre="FIB",Ingres
      =2011)
2 > l[["Centre"]]
3 [1] "FIB"
4 >
```

Com a conseqüència, tenim tres maneres d'accedir al mateix component:

```
1 > l<-list(Nom="Maria Perez",Centre="FIB",Ingres
      =2011)
2 >
3 > l[["Centre"]]
4 [1] "FIB"
5 >
6 > l$Centre
7 [1] "FIB"
8 >
9 > l[[2]]
10 [1] "FIB"
11 >
```

3.3 Afegir i eliminar elements d'un list

Els lists poden créixer i decreïxer dinàmicament. És a dir, podem afegir nous components a un list i eliminar qualsevol dels components que existeixin.

Per a afegir un nou element directament li assignem un valor al component i això ja crea el component:

```
1 > l<-list(Nom="Maria Perez",Centre="FIB",Ingres=2011)
2 > l
3 $Nom
4 [1] "Maria Perez"
5
6 $Centre
7 [1] "FIB"
8
9 $Ingres
10 [1] 2011
11
12 # Afegim l'element "Universitat":
13
14 > l$Universitat <- "UPC"
15 > l
16 $Nom
17 [1] "Maria Perez"
18
19 $Centre
20 [1] "FIB"
21
22 $Ingres
23 [1] 2011
24
25 $Universitat
26 [1] "UPC"
27 >
```

Per a eliminar un element el que farem és assignar-li el valor NULL al component que volem eliminar:

```
1 > l<-list(Nom="Maria Perez",Centre="FIB",Ingres=2011,
2           Universitat="UPC")
3 > l
4 $Nom
5 [1] "Maria Perez"
6
7 $Centre
8 [1] "FIB"
9
10 $Ingres
11 [1] 2011
12
13 $Universitat
```

```

13 [1] "UPC"
14
15 # Eliminem l'element "Centre":
16
17 > l$Centre<-NULL
18 > l
19 $Nom
20 [1] "Maria Perez"
21
22 $Ingres
23 [1] 2011
24
25 $Universitat
26 [1] "UPC"
27 >

```

També es pot preguntar si un element existeix en un list o no mitjançant la funció *is.null*:

is.null : *component_del_list* → *bool*

```

1 > is.null(l$Centre)
2 [1] TRUE
3
4 > is.null(l[["Ingres"]])
5 [1] FALSE
6 >

```

3.4 Accés a sublists

Per obtenir un sublist a partir d'un list, es fa servir la següent forma d'indexada:

$[] : list \times interval \rightarrow list$

on *interval* és un interval de la forma $[exp]$ o $[exp1 : exp2]$. La crida es fa com segueix:

$l[exp]$

```

1 > l<-list(Nom="Maria Perez",Centre="FIB",Ingres =2011,
2         Universitat="UPC")
3 > l[1]
4 $Nom
5 [1] "Maria Perez"
6 >
7 > l[[1]]
8 [1] "Maria Perez"
9 >

```

S'ha de notar la diferència de significat quan es fa servir gafet simple i quan es fa servir doble gafet. En el primer cas torna un list mentre que en el segon cas torna l'èsim component.

També podem fer un sublist amb més d'un element del list original usant un subrang per a l'índex:

```
l[exp1:exp2]
```

```

1 > l[1:3]
2 $Nom
3 [1] "Maria Perez"
4
5 $Centre
6 [1] "FIB"
7
8 $Ingres
9 [1] 2011
10
11 > l[3:4]
12 $Ingres
13 [1] 2011
14
15 $Universitat
16 [1] "UPC"
17 >
18 > l[3:1]
19 $Ingres
20 [1] 2011
21
22 $Centre
23 [1] "FIB"
24
25 $Nom
26 [1] "Maria Perez"
27 >
28 > l[(2+1):(3+1)]
29 $Ingres
30 [1] 2011
31
32 $Universitat
33 [1] "UPC"
34 >

```

3.5 Recorreguts i cerques

Hi ha molts problemes on s'han de recórrer els components d'un list i també problemes on és necessari fer servir l'esquema de cerca. Aquestes dues famílies de problemes les hem caracteritzat al Capítol 1, quan parlàvem de vectors, a les seccions 1.3 i 1.4. Com a exemple, aquí podem veure una funció que escriu un list:

```
1 escriure_list <- function(l){
2   n <- names(l)
3   for (i in 1:length(l)){
4     cat(n[i],": ",l[[i]],"\n")
5   }
6 }
```

que fa un recorregut del list i que es pot utilitzar d'aquesta forma:

```
1 > l <- list(Nom="Joan",Cognom="Martinez",Edat=23,Grau=
   "Informatica",Universitat="UPC")
2 > escriure_list(l)
3 Nom : Joan
4 Cognom : Martinez
5 Edat : 23
6 Grau : Informatica
7 Universitat : UPC
8 >
```

Capítol 4

Data Frames

Un **data frame** df és un objecte **estructurat**, de dues dimensions, on les columnes (tambè anomenades *variables*) són dades d'un mateix tipus i les files (o *individus*) són un conjunt de dades de diferent tipus. Un data frame es pot veure com un list de vectors i també com una matriu no heterogènia. Vist com una matriu, un data frame és un objecte compost per un nombre finit de files F (entrades horitzontals) i un nombre finit de columnes C (entrades verticals) anomenats components o elements. De manera genèrica, podem considerar que un data frame df és de la forma:

$$\begin{array}{c} \overbrace{\hspace{10em}}^{C \text{ columnes}} \\ \left[\begin{array}{ccccc} df_{1,1} & df_{1,2} & \dots & df_{1,C-1} & df_{1,C} \\ df_{2,1} & df_{2,2} & \dots & df_{2,C-1} & df_{2,C} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ df_{F-1,1} & df_{F-1,2} & \dots & df_{F-1,C-1} & df_{F-1,C} \\ df_{F,1} & df_{F,2} & \dots & df_{F,C-1} & df_{F,C} \end{array} \right] \\ \begin{array}{c} F \text{ files} \end{array} \end{array}$$

Com podem veure, en aquest cas el data frame df té F files i C columnes i per tant té $F \times C$ components on es compleix

$$\forall i \in \{1, \dots, F\}, \forall j \in \{1, \dots, C\}, df_{i,j} \in T_j$$

Dit d'una altra manera, amb el data frame df hom pot associar la dimensió $F \times C$. En un data frame, cada fila i és un objecte de tipus *list* amb dades no necessàriament homogènies mentre que cada columna j és un objecte de tipus vector i, per tant, homogeni. És a dir, $df_{i,j}$ es un valor del tipus base T_j . Per fixar la idea, si considerem una enquesta, un data frame és el tipus d'objecte de l'R que permet representar les dades recollides mitjançant l'enquesta. Vist d'aquesta manera, cada columna correspon als valors d'una variable estadística (i possiblement el seu nom és a la capçalera) i cada fila correspon a una resposta donada per una persona enquestada.

Considerem per exemple el data frame següent:

$$df = \begin{bmatrix} Edat & Sexe & Alçada & Hobby \\ 18 & H & 1.81 & Fútbol \\ 20 & M & 1.67 & Lectura \\ 15 & M & 1.60 & Dança \\ 19 & H & 1.76 & Bowling \\ 23 & M & 1.70 & Música \\ 18 & H & 1.65 & Lectura \end{bmatrix}$$

Interpretant aquest exemple com les respostes a una enquesta, les variables estadístiques serien *Edat*, *Sexe*, *Alçada* i *Hobby* on els seus valors serien els que es troben a la columna que encapçala cadascuna. Addicionalment, cada fila seria una possible resposta a l'enquesta. Com hem dit abans, cada fila pot ser vista com un *list*, com ara (23, M, 1.70, Música) i cada columna com a un vector, per exemple el vector de seqüència de caràcters ("Fútbol", L-Lectura", "Dança", "Bowling", "Música", L-Lectura").

4.1 Creació de data frames al llenguatge R

El llenguatge R ofereix diverses formes de crear data frames. Ací en veurem dos: la creació directa i la creació a partir de fitxers.

4.1.1 Creació directa

La creació directa d'un data frame utilitza el constructor

```
data.frame()
```

Aquest constructor rep com a paràmetres les columnes del data frame a crear. Com en el cas dels *list*, hom pot assignar un nom a cadascuna de les columnes (variables) del data frame. Per exemple podem construir un data frame fent:

```
data.frame(nomcol_1 = valcol_1, ..., nomcol_k = valcol_k, stringsAsFactors=FALSE)
```

Un exemple concret de construcció d'un data frames és

```
1 > df <- data.frame(Nom=c("Maria","Kim","Pep"),
2                   Edat=c(18,20,23),
3                   stringsAsFactors=FALSE)
4 > df
5     Nom Edat
6 1  Maria  18
7 2   Kim  20
8 3   Pep  23
9 >
```

Però a diferència del *list*, en aquest cas també poden usar-se els noms nom de les variables d'R per a donar nom als camps del data frame. Així doncs podem fer el següent:

```
1 > Edat <- c(18,20,15,19,23,18)
2 > Sexe <- c("H","M","M","H","M","H")
3 > Alcada <- c(1.81,1.67,1.60,1.76,1.70,1.65)
```



```

4 > Hobby <- c("Futbol","Lectura","Danca","Bowling","
      Musica","Lectura")
5 > df <- data.frame(Edat,Sexe,Alcada,Hobby,
      stringsAsFactors=FALSE)
6 > df
7   Edat Sexe Alcada   Hobby
8 1   18   H   1.81 Futbol
9 2   20   M   1.67 Lectura
10 3   15   M   1.60 Danca
11 4   19   H   1.76 Bowling
12 5   23   M   1.70 Musica
13 6   18   H   1.65 Lectura
14 >

```

Recordem que el paràmetre `stringsAsFactors` defineix la manera en la que R interpretarà les cadenes de caràcters. Si no es defineix de manera explícita el paràmetre `stringsAsFactors=FALSE`, el valor que pren per defecte és `TRUE` i això vol dir que les cadenes de caràcters seran tractades d'una forma diferent a la que estem acostumats. L'explicació en detall d'aquest paràmetre està fora de l'abast d'aquest document i per tant, sempre que calgui crear data frames amb columnes on hi hagi cadenes de caràcters s'inclourà l'assignació `stringsAsFactors=FALSE`.

També es poden crear data frames buits. Per exemple, si a l'intendent R introduïm l'assignació

```

1 > df1 <- data.frame()

```

i tot seguit demanem pel valor de la variable `df1`, s'obté

```

1 > df1
2 data frame with 0 columns and 0 rows
3 >

```

4.1.2 Creació de data frames a partir de fitxers

Generalment els data frames són objectes grans i interessarà construir-los a partir de les dades d'un fitxer CSV, format en el qual les dades estan usualment separades per un espai. En aquest cas, R ofereix un operador de lectura de fitxers `read.table` que s'utilitza com segueix:

```
read.table(quote(nomfitxer), header=TRUE)
```

El paràmetre `header = TRUE` defineix que la primera línia conté les capçaleres del data frame, és a dir, els noms de les variables. Si no s'explicita pren el valor per defecte `FALSE`. Per exemple, el codi R

```

1 > df <- read.table("mydata.txt", header=TRUE)

```

crea el data frame `df` llegint les dades del fitxer de nom `"mydata.txt"` amb capçaleres.

4.2 Operadors del tipus *list* i data frames

Ja hem explicat que un data frame es pot veure com un *list*, per tant se li podran aplicar tots els operadors estudiats al capítol anterior.

4.2.1 Accés a una columna

En tot el que segueix, la variable `df` denotarà un data frame. suposem el data frame de la Secció 4.1.1, la primera columna corresponent al nom `Edat` es pot seleccionar usant el valor enter de l'índex

```
1 > df[[1]]
2 [1] 18 20 15 19 23 18
```

El nom de de la columna amb l'operador `$`

```
1 > df$Edat
2 [1] 18 20 15 19 23 18
```

O bé emprant el nom com a índex

```
1 > df[["Hobby"]]
2 [1] "Futbol" "Lectura" "Danca" "Bowling" "Musica"
   "Lectura"
```

4.2.2 Obtenir la longitud

Podem demanar la longitud del data frame amb la funció `length()` i ens retornarà el nombre de columnes (variables) del data frame:

```
1 > length(df)
2 [1] 4
3 >
```

4.2.3 Noms

Podem accedir als noms de les variables:

names(): Retorna un vector amb els noms de les variables o columnes.

colnames(): Retorna un vector amb els noms de les variables. És equivalent a `names()`.

o també als noms de les files:

rownames(): Retorna un vector amb els noms de les files.

En tots els casos, igual que passa en les estructures *list*, també es poden assignar a aquestes funcions el vector de noms que hom vulgui definir. Exemples d'utilització d'aquestes operacions aplicades al data frame definit al llistat 4.1.1 són:

```
1 > names(df)
2 [1] "Edat" "Sexe" "Alcada" "Hobby"
3 > colnames(df)
4 [1] "Edat" "Sexe" "Alcada" "Hobby"
5 > rownames(df)
6 [1] "1" "2" "3" "4" "5" "6"
7 >
```

Considerem el data frame del llistat de la Secció 4.1.1. Si ara hom defineix el conjunt de noms de files amb

```
1 > rownames(df) <- c("Joan", "Anna", "Maria", "Albert", "
2   Jana", "Josep")
3 >
```

el data frame `df` resultant serà

```
1 > df
2   Edat Sexe Alcada Hobby
3 Joan   18   H   1.81 Futbol
4 Anna   20   M   1.67 Lectura
5 Maria  15   M   1.60 Danca
6 Albert 19   H   1.76 Bowling
7 Jana   23   M   1.70 Musica
8 Josep  18   H   1.65 Lectura
9 > length(df)
10 [1] 4
11 >
```

4.3 Operadors del tipus matriu i data frames

Com ha estat explicat més amunt, un data frame també es pot veure com una matriu, per tant també es podem aplicar els operadors definits per a les matriu i accedir als elements del data frame de forma matricial, és a dir, variables amb dos índexs. Recordem-ne algunes d'aquestes operacions.

nrow(): Retorna el nombre de files.

ncol(): Retorna el nombre de columnbes.

dim(): Retorna el nombre de files i columnbes, és a dir, les dimensions de la matriu.

Considerem ara el data frame `df` del llistat de la Secció 4.2.3. Aleshores tindrem

```
1 > nrow(df)
2 [1] 6
3 > ncol(df)
4 [1] 4
5 > dim(df)
6 [1] 6 4
7 >
```

Important En particular, quan s'usen operadors d'indexat, cal assegurar-se'n que el valor resultant de l'avaluació de l'expressió de l'índex caigui dins d'un rang de valors `[1..max]` on `max` és el valor màxim de files o columnbes del data frame. En general cal tenir molta cura amb el fet que les funcions s'apliquin sobre data frames que estiguin correctament definits. Un exemple concret de situació en la qual una funció s'aplica de manera errònia és el següent

```

1 > df2
2 data frame with 0 columns and 0 rows
3 > nrow(df2)
4 [1] 0
5 > ncol(df2)
6 [1] 0
7 > df[1,2]
8 NULL

```

4.3.1 Accés a una columna sencera

Recordem que cada columna d'un data frame és un vector, per tant, un objecte amb informació homogènia. Considerem un altre cop el data frame `df` del llistat de la Secció 4.2.3. Aleshores tindrem

```

1 > df
2      Edat Sexe Alcada  Hobby
3 Joan   18   H   1.81  Futbol
4 Anna  20   M   1.67  Lectura
5 Maria 15   M   1.60   Danca
6 Albert 19   H   1.76  Bowling
7 Jana  23   M   1.70  Musica
8 Josep 18   H   1.65  Lectura
9 > df[,2]
10 [1] "H" "M" "M" "H" "M" "H"
11 > df[,3]
12 [1] 1.81 1.67 1.60 1.76 1.70 1.65
13 >

```

4.3.2 Accés a una fila sencera

Recordem que cada fila és un list per tant, un objecte heterogeni. Fixeu-vos que en fer la petició d'una única fila en realitat ens retorna un data frame amb només aquesta fila.

```

1 > df
2      Edat Sexe Alcada  Hobby
3 Joan   18   H   1.81  Futbol
4 Anna  20   M   1.67  Lectura
5 Maria 15   M   1.60   Danca
6 Albert 19   H   1.76  Bowling
7 Jana  23   M   1.70  Musica
8 Josep 18   H   1.65  Lectura
9 > df[2,]
10      Edat Sexe Alcada  Hobby
11 Anna  20   M   1.67  Lectura
12 > df[4,]
13      Edat Sexe Alcada  Hobby
14 Albert 19   H   1.76  Bowling

```

4.3.3 Accés a una component d'un data frame

De manera anàloga al que passava en el càlcul matricial, tot sovint interessa tractar cada element d'un data frame de manera individual. Considerant el data frame `df` del llistat de la Secció 4.2.3 tindrem

```

1 > df
2       Edat Sexe Alcada  Hobby
3 Joan    18   H   1.81  Futbol
4 Anna    20   M   1.67  Lectura
5 Maria   15   M   1.60   Danca
6 Albert  19   H   1.76  Bowling
7 Jana    23   M   1.70  Musica
8 Josep   18   H   1.65  Lectura
9 > df[1,3]
10 [1] 1.81
11 > df[3,4]
12 [1] "Danca"
13 > df[2,5]
14 NULL
15 > df[7,3]
16 NULL
17 >

```

Noteu que, quan es tracta d'accedir a elements del data frame que no estan definits, l'operador retorna la constant `NULL`, és a dir, un valor invàlid. Recordeu el que ha estat dit més amunt respecte d'haver d'usar amb cura els índexs.

També es pot accedir a un element determinat usant com a índex el nom de la columna

```

1 > df
2       Edat Sexe Alcada  Hobby
3 Joan    18   H   1.81  Futbol
4 Anna    20   M   1.67  Lectura
5 Maria   15   M   1.60   Danca
6 Albert  19   H   1.76  Bowling
7 Jana    23   M   1.70  Musica
8 Josep   18   H   1.65  Lectura
9 > df[1,"Alcada"]
10 [1] 1.81
11 > df[3,"Hobby"]
12 [1] "Danca"
13 >

```

4.4 Subdata frames

De la mateixa manera que de les matrius es poden extreure submatrius, dels data frames es poden extreure subdata frames. Les diverses extraccions es poden fer per rang d'índex, amb filtres o selecció de files

4.4.1 Selecció per rang d'índex

La selecció de subdata frames per rang d'índex segueix la mateixa sintaxi que en el cas de l'extracció de submatrius

- L'absència del segon índex, `[i,]` designa la fila *i*-èsima.
- L'absència del primer índex, `[, j]` designa la columna *j*-èsima.
- La presència d'un rang, `[min:max]` indica el conjunt de files o columnes tals que la primera és la corresponent a l'índex `min` i la darrera és la corresponent a l'índex `max`.

Vegem tot seguit alguns exemples.

```

1 > df
2       Edat Sexe Alcada  Hobby
3 Joan    18   H   1.81  Futbol
4 Anna    20   M   1.67  Lectura
5 Maria   15   M   1.60   Danca
6 Albert  19   H   1.76  Bowling
7 Jana    23   M   1.70  Musica
8 Josep   18   H   1.65  Lectura
9 >
10 > df[2:4,]
11       Edat Sexe Alcada  Hobby
12 Anna    20   M   1.67  Lectura
13 Maria   15   M   1.60   Danca
14 Albert  19   H   1.76  Bowling
15 >
16 > df[,1:3]
17       Edat Sexe Alcada
18 Joan    18   H   1.81
19 Anna    20   M   1.67
20 Maria   15   M   1.60
21 Albert  19   H   1.76
22 Jana    23   M   1.70
23 Josep   18   H   1.65
24 >
25 > df[2:4,1:3]
26       Edat Sexe Alcada
27 Anna    20   M   1.67
28 Maria   15   M   1.60
29 Albert  19   H   1.76
30 >
31 > df[2:4,3]
32 [1] 1.67 1.60 1.76

```

Noteu que el resultat de la darrera línia és un vector.

Com hem vist, l'accés a les dades d'una única columna del data frame retorna les dades de la columna en un vector. Si el que volem és que el resultat retornat sigui una columna del data frame expressada com un altre data frame cal afegir el paràmetre `drop = FALSE` com segueix:

```

1 > df[2:4,3, drop=FALSE]
2     Alcada
3 Anna    1.67
4 Maria   1.60
5 Albert  1.76

```

4.4.2 Filtres

Moltes vegades el que es vol és obtenir un subdata frame segons alguna condició o propietat que es pot expressar com una expressió lògica respecte d'una columna. Si el data frame és `df`, la sentència seria:

```
df[condició(df$colnom, valor), rangcol]
```

per exemple

```

1 > df
2     Edat Sexe Alcada  Hobby
3 Joan   18   H   1.81  Futbol
4 Anna   20   M   1.67  Lectura
5 Maria  15   M   1.60   Danca
6 Albert 19   H   1.76  Bowling
7 Jana   23   M   1.70  Musica
8 Josep  18   H   1.65  Lectura
9 > df[df$Alcada < 1.70,]
10     Edat Sexe Alcada  Hobby
11 Anna   20   M   1.67  Lectura
12 Maria  15   M   1.60   Danca
13 Josep  18   H   1.65  Lectura
14 > df[df$Alcada > 1.70, 2:4]
15     Sexe Alcada  Hobby
16 Joan   H   1.81  Futbol
17 Albert H   1.76  Bowling
18 >

```

4.4.3 Selecció de files

Si hom vol seleccionar un subdata frame que contingui totes les columnes del data frame original però només inclogui les files que compleixin una determinada condició, el llenguatge R ofereix la funció `subset`. La sintaxi és

```
subset(df, cond(colnoms, valors))
```

Exemples d'utilització són

```

1 > df
2     Edat Sexe Alcada  Hobby
3 Joan   18   H   1.81  Futbol
4 Anna   20   M   1.67  Lectura
5 Maria  15   M   1.60   Danca
6 Albert 19   H   1.76  Bowling
7 Jana   23   M   1.70  Musica
8 Josep  18   H   1.65  Lectura
9 >

```

```

10 > subset(df, Sexe == "H")
11       Edat Sexe Alcada  Hobby
12 Joan    18   H   1.81  Futbol
13 Albert  19   H   1.76 Bowling
14 Josep   18   H   1.65 Lectura
15 > subset(df, Sexe == "H" & Alcada >= 1.65)
16       Edat Sexe Alcada  Hobby
17 Joan    18   H   1.81  Futbol
18 Albert  19   H   1.76 Bowling
19 Josep   18   H   1.65 Lectura
20 >

```

4.5 Extensió de data frames

Els data frame es poden estendre afegint tant files com columnes usant els operadors específics del llenguatge R.

4.5.1 Afegir una fila

Afegir una fila a un data frame significa afegir un individu, és a dir un valor concret per a cadascuna de les variables. La funció corresponent és

```
rbind()
```

Per a poder afegir un individu (una fila) al data frame, com que les dades a afegir són heterogènies, possiblement de diferents tipus, cal afegir un altre *data frame* o un *list*. Aquest nou element cal que tingui necessàriament el mateix nombre de columnes que el data frame que es vol augmentar. Un exemple que il·lustra aquesta funció és:

```

1 > df
2       Edat Sexe Alcada  Hobby
3 Joan    18   H   1.81  Futbol
4 Anna    20   M   1.67 Lectura
5 Maria   15   M   1.60  Danca
6 Albert  19   H   1.76 Bowling
7 Jana    23   M   1.70 Musica
8 Josep   18   H   1.65 Lectura
9 > rbind(df, data.frame(Edat=22, Sexe="M", Alcada=1.70,
10       Hobby="Tenis",
11       stringsAsFactors=FALSE))
12       Edat Sexe Alcada  Hobby
13 Joan    18   H   1.81  Futbol
14 Anna    20   M   1.67 Lectura
15 Maria   15   M   1.60  Danca
16 Albert  19   H   1.76 Bowling
17 Jana    23   M   1.70 Musica
18 Josep   18   H   1.65 Lectura
19 7       22   M   1.70  Tenis
20 >

```


Fixeu-vos que, afegint data frame, cal tenir cura de definir l'opció `stringsAsFactors=FALSE` si algun dels camps del data frame és una cadena de caràcters.

Afegir una nova fila a un data frame expressada com un *list* s'aconsegueix amb la sintaxi

```

1 > df
2       Edat Sexe Alcada  Hobby
3 Joan   18   H   1.81  Futbol
4 Anna   20   M   1.67  Lectura
5 Maria  15   M   1.60   Danca
6 Albert 19   H   1.76  Bowling
7 Jana   23   M   1.70  Musica
8 Josep  18   H   1.65  Lectura
9 > NovaFila <- list(Edat=22, Sexe="M", Alcada=1.70,
10                   Hobby="Tenis")
11 > rbind(df, NovaFila)
12       Edat Sexe Alcada  Hobby
13 Joan   18   H   1.81  Futbol
14 Anna   20   M   1.67  Lectura
15 Maria  15   M   1.60   Danca
16 Albert 19   H   1.76  Bowling
17 Jana   23   M   1.70  Musica
18 Josep  18   H   1.65  Lectura
19 7       22   M   1.70   Tenis

```

També es poden afegir diverses files, és a dir, es pot afegir un data frame a un altre data frame.

```

1 > df
2       Edat Sexe Alcada  Hobby
3 Joan   18   H   1.81  Futbol
4 Anna   20   M   1.67  Lectura
5 Maria  15   M   1.60   Danca
6 Albert 19   H   1.76  Bowling
7 Jana   23   M   1.70  Musica
8 Josep  18   H   1.65  Lectura
9 >
10 > df1 <- data.frame(Edat=c(17,21,15), Sexe=c("H","M","H"),
11                   Alcada=c(1.70,1.71,1.65),
12                   Hobby=c("Bowling","Danca",
13                           "Musica"),
14                   stringsAsFactors=FALSE)
15 >
16 > df1
17       Edat Sexe Alcada  Hobby
18 1       17   H   1.70  Bowling
19 2       21   M   1.71   Danca
20 3       15   H   1.65  Musica

```

```

20 > rbind(df, df1)
21      Edat Sexe Alcada  Hobby
22 Joan   18   H   1.81  Futbol
23 Anna   20   M   1.67  Lectura
24 Maria  15   M   1.60   Danca
25 Albert 19   H   1.76  Bowling
26 Jana   23   M   1.70  Musica
27 Josep  18   H   1.65  Lectura
28 7      17   H   1.70  Bowling
29 8      21   M   1.71   Danca
30 9      15   H   1.65  Musica
31 >

```

4.5.2 Afegir una nova columna

Afegir una columna a un data frame significa afegir una variable, és a dir, afegir un nou valor a tots i cadascun dels individus del data frame. La funció és

```
cbind()
```

Per a poder afegir una variable o columna al data frame, les dades poden estar emmagatzemades directament en un vector, perquè són homogènies. L'única restricció que cal tenir en compte és que el nombre de components del vector ha de ser igual al nombre de files del data frame.

```

1 > Ciutat <- c("Paris", "Barcelona", "Barcelona", "Roma", "
      Caracas", "Barcelona", "Paris")
2 > df <- rbind(df, list(22, "M", 1.70, "Tenis"))
3 > df <- cbind(df, Ciutat)
4 >
5 > df
6      Edat Sexe Alcada  Hobby  Ciutat
7 Joan   18   H   1.81  Futbol   Paris
8 Anna   20   M   1.67  Lectura Barcelona
9 Maria  15   M   1.60   Danca Barcelona
10 Albert 19   H   1.76  Bowling   Roma
11 Jana   23   M   1.70  Musica   Caracas
12 Josep  18   H   1.65  Lectura Barcelona
13 7      22   M   1.70   Tenis   Paris
14 >

```

També es pot afegir una columna al data frame tal i com es feia en el cas del *list*, és a dir, aplicant directament l'operador "\$":

```

1 > df
2      Edat Sexe Alcada  Hobby
3 Joan   18   H   1.81  Futbol
4 Anna   20   M   1.67  Lectura
5 Maria  15   M   1.60   Danca
6 Albert 19   H   1.76  Bowling
7 Jana   23   M   1.70  Musica
8 Josep  18   H   1.65  Lectura

```

```

9 > df$CodiPostal <- c
      (08027,08003,08014,08034,08034,08006)
10 > df
11      Edat Sexe Alcada Hobby CodiPostal
12 Joan    18   H   1.81  Futbol      8027
13 Anna    20   M   1.67 Lectura      8003
14 Maria   15   M   1.60  Danca       8014
15 Albert  19   H   1.76 Bowling     8034
16 Jana    23   M   1.70 Musica      8034
17 Josep   18   H   1.65 Lectura     8006
18 >

```

Notis que, a diferència de la funció `cbind()`, afegir directament columnes modifica el data frame.

4.6 Modificació de data frames

Hi ha alguns problemes on cal modificar parcialment els data frames, és a dir, cal modificar els valors d'algunes variables. En aquest cas, la funció de l'R que permet fer aquestes modificacions és

```
transform()
```

A continuació podem veure alguns exemples d'ús d'aquesta funció.

```

1 > df
2      Edat Sexe Alcada Hobby CodiPostal
3 Joan    18   H   1.81  Futbol      8027
4 Anna    20   M   1.67 Lectura      8003
5 Maria   15   M   1.60  Danca       8014
6 Albert  19   H   1.76 Bowling     8034
7 Jana    23   M   1.70 Musica      8034
8 Josep   18   H   1.65 Lectura     8006
9 > transform(df, Alcada=Alcada*100)
10      Edat Sexe Alcada Hobby CodiPostal
11 Joan    18   H   181  Futbol      8027
12 Anna    20   M   167 Lectura      8003
13 Maria   15   M   160  Danca       8014
14 Albert  19   H   176 Bowling     8034
15 Jana    23   M   170 Musica      8034
16 Josep   18   H   165 Lectura     8006
17 > transform(df,Alcada=Alcada%%2.54) #passem cm a
      polzades
18      Edat Sexe Alcada Hobby CodiPostal
19 Joan    18   H   71.26  Futbol      8027
20 Anna    20   M   65.75 Lectura      8003
21 Maria   15   M   62.99  Danca       8014
22 Albert  19   H   69.29 Bowling     8034
23 Jana    23   M   66.93 Musica      8034
24 Josep   18   H   64.96 Lectura     8006
25 > transform(df,Alcada=Alcada%%12,CodiPostal=
      CodiPostal%%100) #passem polzades a peus, reduim cp

```

```

26           Edat Sexe Alcada   Hobby CodiPostal
27 Joan      18   H   5.94   Futbol      27
28 Anna     20   M   5.48  Lectura      3
29 Maria    15   M   5.25   Danca      14
30 Albert   19   H   5.77  Bowling     34
31 Jana     23   M   5.58   Musica     34
32 Josep    18   H   5.41  Lectura      6
33 >

```

4.7 Ordenació de data frames

Els data frames poden ser ordenats segons una o diverses columnes. La funció que cal utilitzar és

```
order()
```

L'ordenació es fa de manera creixent o decreixent en funció del paràmetre `decreasing` que pren valor cert o fals. Per defecte pren valor `FALSE` i l'ordenació és fa segons valors creixents. Per exemple, hom pot ordenar el data frame `df` de més amunt segons la variable `Edat` i en ordre creixent com segueix

```

1 > df
2           Edat Sexe Alcada   Hobby CodiPostal
3 Joan      18   H   1.81   Futbol     8027
4 Anna     20   M   1.67  Lectura     8003
5 Maria    15   M   1.60   Danca     8014
6 Albert   19   H   1.76  Bowling     8034
7 Jana     23   M   1.70   Musica     8034
8 Josep    18   H   1.65  Lectura     8006
9 >
10 > df[order(df$Edat),]
11           Edat Sexe Alcada   Hobby CodiPostal
12 Maria    15   M   1.60   Danca     8014
13 Joan     18   H   1.81   Futbol     8027
14 Josep    18   H   1.65  Lectura     8006
15 Albert   19   H   1.76  Bowling     8034
16 Anna     20   M   1.67  Lectura     8003
17 Jana     23   M   1.70   Musica     8034
18 >

```

Notis la coma que apareix al darrere de la columna seleccionada `$Edat)`,. L'ordenació decreixent seria

```

1 > df
2           Edat Sexe Alcada   Hobby CodiPostal
3 Joan      18   H   1.81   Futbol     8027
4 Anna     20   M   1.67  Lectura     8003
5 Maria    15   M   1.60   Danca     8014
6 Albert   19   H   1.76  Bowling     8034
7 Jana     23   M   1.70   Musica     8034
8 Josep    18   H   1.65  Lectura     8006
9 >

```

```

10 > df[order(df$Edat, decreasing=TRUE),]
11      Edat Sexe Alcada  Hobby CodiPostal
12 Jana    23    M   1.70  Musica      8034
13 Anna    20    M   1.67  Lectura      8003
14 Albert  19    H   1.76  Bowling      8034
15 Joan    18    H   1.81  Futbol      8027
16 Josep   18    H   1.65  Lectura      8006
17 Maria   15    M   1.60   Danca      8014
18 >

```

Una ordenació creixent simultàniament segons Edat i Alcada seria

```

1 > df
2      Edat Sexe Alcada  Hobby CodiPostal
3 Joan    18    H   1.81  Futbol      8027
4 Anna    20    M   1.67  Lectura      8003
5 Maria   15    M   1.60   Danca      8014
6 Albert  19    H   1.76  Bowling      8034
7 Jana    23    M   1.70  Musica      8034
8 Josep   18    H   1.65  Lectura      8006
9 > df[order(df$Edat, df$Alcada),]
10      Edat Sexe Alcada  Hobby CodiPostal
11 Maria   15    M   1.60   Danca      8014
12 Josep   18    H   1.65  Lectura      8006
13 Joan    18    H   1.81  Futbol      8027
14 Albert  19    H   1.76  Bowling      8034
15 Anna    20    M   1.67  Lectura      8003
16 Jana    23    M   1.70  Musica      8034
17 >

```

També podríem afegir, ahora, una selecció per columnes i mostrar només un subrang d'elles, com per exemple:

```

1 > df
2      Edat Sexe Alcada  Hobby CodiPostal
3 Joan    18    H   1.81  Futbol      8027
4 Anna    20    M   1.67  Lectura      8003
5 Maria   15    M   1.60   Danca      8014
6 Albert  19    H   1.76  Bowling      8034
7 Jana    23    M   1.70  Musica      8034
8 Josep   18    H   1.65  Lectura      8006
9 > df[order(df$Edat, df$Alcada), 2:4]
10      Sexe Alcada  Hobby
11 Maria    M   1.60   Danca
12 Josep    H   1.65  Lectura
13 Joan     H   1.81  Futbol
14 Albert   H   1.76  Bowling
15 Anna     M   1.67  Lectura
16 Jana     M   1.70  Musica
17 >

```

4.8 Fusió de data frames

El llenguatge R permet construir data frames a partir de la fusió de dos data frames donats. La fusió es fa a partir d'una variable (columna) que els data frames tinguin en comú. El data frame resultant estarà compost per

1. La unió de les columnes dels dos data frames. Notis que això vol dir que no es repeteixen les que tinguin en comú.
2. Les files dels dos data frames tals que tenen el mateix índex de fila i el valor de la columna que tenen en comú els data frames són iguals.

La sintaxi de la funció és

```
merge(df1, df2)
```

Un exemple de data frame resultant de la fusió dels data frames `df1` i `df2` és

```

1 > df1
2   Var1 Var2 Var3
3 1   A1   B4   C1
4 2   A2   B2   C2
5 3   A3   B1   C3
6 4   A4   B2   C1
7 5   A2   B7   C3
8 6   A4   B3   C4
9 7   A5   B6   C3
10 8   A6   B1   C3
11 >
12 > df2
13   Var2 Var4
14 1   B1   D1
15 2   B2   D2
16 3   B3   D1
17 4   B4   D3
18 5   B5   D5
19 >
20 >
21 > merge(df1, df2)
22 >
23   Var2 Var1 Var3 Var4
24 1   B1   A3   C3   D1
25 2   B1   A6   C3   D1
26 3   B2   A2   C2   D2
27 4   B2   A4   C1   D2
28 5   B3   A4   C4   D1
29 6   B4   A1   C1   D3
30 >
```

També podem fer servir el paràmetre `all` per indicar que el *merge* inclogui totes les files malgrat que hi hagi valors de la columna compartida que no siguin iguals. Per exemple

```

1 > df1
2   Var1 Var2 Var3
3 1   A1   B4   C1
4 2   A2   B2   C2
5 3   A3   B1   C3
6 4   A4   B2   C1
7 5   A2   B7   C3
8 6   A4   B3   C4
9 7   A5   B6   C3
10 8   A6   B1   C3
11 >
12 > df2
13 >
14   Var2 Var4
15 1   B1   D1
16 2   B2   D2
17 3   B3   D1
18 4   B4   D3
19 5   B5   D5
20 >
21 > merge(df1, df2, all = TRUE)
22 >
23   Var2 Var1 Var3 Var4
24 1   B1   A3   C3   D1
25 2   B1   A6   C3   D1
26 3   B2   A2   C2   D2
27 4   B2   A4   C1   D2
28 5   B3   A4   C4   D1
29 6   B4   A1   C1   D3
30 7   B5 <NA> <NA>   D5
31 8   B6   A5   C3 <NA>
32 9   B7   A2   C3 <NA>
33 >

```

Notis que quan a les files on hi ha valors no compartits en la columna `Var2` prenen per valor `<NA>`.

Quan hi ha columnes (variables) que tenen valors comuns a dos data frames però els seus noms no coincideixin es pot aplicar l'operador `merge` fent servir els paràmetres `by.x` i `by.y` per indicar el nom de la variable a considerar tan a l'operador `x` com a l'operador `y`.

```

1 > df1
2   Var1 Var2 Var3
3 1   A1   B4   C1
4 2   A2   B2   C2
5 3   A3   B1   C3
6 4   A4   B2   C1
7 5   A2   B7   C3
8 6   A4   B3   C4
9 7   A5   B6   C3
10 8   A6   B1   C3

```

```
11 > df2
12 >
13   Dif2 Var4
14 1   B1   D1
15 2   B2   D2
16 3   B3   D1
17 4   B4   D3
18 5   B5   D5
19 >
20 > merge(df1, df2, by.x = "Var2", by.y = "Dif2")
21 >
22   Var2 Var1 Var3 Var4
23 1   B1   A3   C3   D1
24 2   B1   A6   C3   D1
25 3   B2   A2   C2   D2
26 4   B2   A4   C1   D2
27 5   B3   A4   C4   D1
28 6   B4   A1   C1   D3
29 >
```


Bibliografia

- [1] Wirth, Niklaus.: Algorithms and data structures. London et al.: Prentice-Hall, 1986.