

## 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:

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

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  $\langle 23, M, 1.70, Música \rangle$  i cada columna com a un vector, per exemple el vector de seqüència de caràcters  $\langle "Fútbol", L-Lectura", "Dança", "Bowling", "Música", L-Lectura \rangle$ .

## 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'interpret 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", "
    Jana", "Josep")
2 >

```

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 columnes, 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 columnes 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   Tennis

```

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 Tennis 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, alhora, 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 >
```