

## Capítol 4

# Data Frames

Un **data frame**  $df$  és un objecte **estructurat**, de dues dimensions, on les columnes són dades d'un mateix tipus (de les columnes se'n diuen *variables*) i les files són un conjunt de dades de diferent tipus (de les files se'n diuen *individus*). 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:

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

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  $\forall i \in \{1, \dots, F\} \forall j \in \{1, \dots, C\}, df_{i,j} \in T_j$ . Dit d'altra manera, al data frame  $df$  li podem associar la dimensió  $F \times C$ . En un data frame, cada fila  $i$  és un objecte de tipus list (no homogeni) mentre que cada columna  $j$  és un objecte de tipus vector (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 ens permet representar les dades recollides mitjançant l'aplicació de 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:

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

Veient aquest exemple com 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 la 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  $\langle Fútbol, Lectura, Dança, Bowling, Música, Lectura \rangle$ .

## 4.1 Creació de data frames al llenguatge R

Per a construir un data frame, usarem l'operador següent:

1. `data.frame()`

Aquest constructor rep com a paràmetres les columnes del data frame a crear. Com en el cas dels list, podem donar nom a aquestes columnes (variables) del data frame, per tant podem construir un data frame fent:

`data.frame(nomcol1 = valcol1, ..., nomcolk = valcolk, stringsAsFactors = FALSE)`

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

Però a diferència del list, en aquest cas també podem aprofitar-nos del nom de les variables de 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 paràmetre *stringsAsFactors* indica a l'R com s'han de tractar les cadenes de caràcters. Si no li posem *stringsAsFactors = FALSE*, per defecte, tindrà el valor *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 necessitem crear data frames amb columnes on hi hagi cadenes de caràcters hem de posar *stringsAsFactors = FALSE*.

També podem crear un data frame buit.

```

1 > df1 <- data.frame()
2 > df1
3 data frame with 0 columns and 0 rows
4 >

```

2. Generalment els data frames són objectes grans i ens interessarà construir-los a partir dels dades d'un fitxer CSV (usualment les dades estan separades per un espai). En aquest cas, fem servir un operador de lectura de fitxers `read.table` que s'utilitza com segueix:

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

El paràmetre *header = TRUE* es fa servir per dir que la primera línia conté les capçaleres del data frame (els noms de les variables). Si no es posa, per defecte es *header = FALSE*.

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

## 4.2 Fent servir operadors del tipus list

Ja hem explicat que un data frame es pot veure com un list, per tant podrem fer servir tots els operadors que hem vist al capítol anterior.

1. `names()`

```

1 > names(df)
2 [1] "Edat" "Sexe" "Alcada" "Hobby"
3 >

```

2. Accés a cada columna:

```

1 > df[[1]]
2 [1] 18 20 15 19 23 18
3 > df$Edat
4 [1] 18 20 15 19 23 18
5 > df[["Hobby"]]
6 [1] "Futbol" "Lectura" "Danca" "Bowling" "
    Musica" "Lectura"
7 >

```

### 4.3 Fent servir operadors del tipus matriu

Però com ja hem dit un data frame també es pot veure com una matriu, per tant també podem fer servir operadors del tipus matriu i accedir als elements de forma matricial.

1. `nrow()`, `ncol()` i `dim()`

```

1 > df
2   Edat Sexe Alcada Hobby
3 1   18   H   1.81 Futbol
4 2   20   M   1.67 Lectura
5 3   15   M   1.60 Danca
6 4   19   H   1.76 Bowling
7 5   23   M   1.70 Musica
8 6   18   H   1.65 Lectura
9 > nrow(df)
10 [1] 6
11 > ncol(df)
12 [1] 4
13 > dim(df)
14 [1] 6 4
15 >

```

**Important** En particular, quan fem servir operadors d'indexat, hem de vigilar que l'avaluació de l'expressió corresponent a l'índex es trobi al rang de valors adequats. En general hem de tenir molta cura amb el fet que les funcions utilitzades tinguin un valor definit per als paràmetres. Seguidament podem veure un exemple de problemes d'aquest tipus.

```

1 > df1
2 data frame with 0 columns and 0 rows
3 > nrow(df1)
4 [1] 0
5 > ncol(df1)
6 [1] 0
7 > ndim(df1)
8 Error: no se pudo encontrar la funcion "ndim"

```

## 2. Accés a una columna sencera.

Recodem que cada columna és un vector, per tant, un objecte homogeni.

```

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

```

## 3. Accés a una fila sencera.

Recodem que cada fila és un list per tant, un objecte heterogeni.

```

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

```

## 4. Accés a un component d'un data frame (accés matricial):

```

1 > df
2   Edat Sexe Alcada   Hobby
3 1   18   H   1.81 Futbol
4 2   20   M   1.67 Lectura
5 3   15   M   1.60  Danca
6 4   19   H   1.76 Bowling
7 5   23   M   1.70 Musica
8 6   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[5,6]
16 NULL
17 >

```

Noteu que quan tractem d'accedir a elements que no estan definits al data frame, l'operador ens torna *NULL*, que no és un objecte vàlid. Això vol dir que s'ha de tenir molta cura quan es fa servir l'indexat.

## 4.4 Subdata frames

1. De la mateixa manera que de les matrius podíem extreure submatrius, dels data frames podem extreure subdata frames. Els operands a fer servir són els mateixos que fèiem servir per extreure submatrius.

```

1 > df
2   Edat Sexe Alcada   Hobby
3 1   18   H   1.81 Futbol
4 2   20   M   1.67 Lectura
5 3   15   M   1.60  Danca
6 4   19   H   1.76 Bowling
7 5   23   M   1.70 Musica
8 6   18   H   1.65 Lectura
9 > df[2:4,]
10  Edat Sexe Alcada   Hobby
11 2   20   M   1.67 Lectura
12 3   15   M   1.60  Danca
13 4   19   H   1.76 Bowling
14 > df[,1:3]
15  Edat Sexe Alcada
16 1   18   H   1.81
17 2   20   M   1.67
18 3   15   M   1.60
19 4   19   H   1.76
20 5   23   M   1.70
21 6   18   H   1.65
22 > df[2:4,1:3]
23  Edat Sexe Alcada
24 2   20   M   1.67
25 3   15   M   1.60
26 4   19   H   1.76
27 > df[2:4,3]
28 [1] 1.67 1.60 1.76 #Noteu que en aquest cas el
    resultat es un vector

```

Com hem vist, l'accés a les dades d'una única columna del data frame ens retorna aquestes dades en un vector. Si el que volem és que ens torni un data frame amb una columna en lloc d'un vector el que cal és afegir el paràmetre *drop = FALSE* en l'accés com segueix:

```

1 > df[2:4,3, drop=FALSE]

```

```

2   Alcada
3 2   1.67
4 3   1.60
5 4   1.76
6 >

```

## 2. Filtres

Moltes vegades el que volem és filtrar els data frames amb una condició lògica respecte d'una columna. És a dir, extreure un subdata frame que satisfaci una condició. Sigui *D* un data frame, la sentència seria:

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

```

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

```

## 3. Selecció de files.

També es pot fer una selecció utilitzant la funció `subset`, que permet fer un subdata frame amb el mateix nombre de columnes però només aquelles files que compleixin la condició:

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

```

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

```

```

14 > subset(df, Sexe == "H" & Alcada >= 1.65)
15   Edat Sexe Alcada   Hobby
16  1   18   H   1.81 Futbol
17  4   19   H   1.76 Bowling
18  6   18   H   1.65 Lectura
19 >

```

## 4.5 Extensió de data frames

Els data frame es poden estendre afegint tant files com columnes fent servir els següents operadors.

1. `rbind()` per afegir una fila, és a dir un individu.

Per a poder afegir un individu (una fila) al data frame, com que les dades a afegir són heterogènies (de diferent tipus), el que s'ha d'afegir és un altre data frame o un list. Aquest nou data frame (o list) ha de tenir necessàriament el mateix nombre de columnes que el data frame que es vol augmentar.

```

1 > df
2   Edat Sexe Alcada   Hobby
3  1   18   H   1.81 Futbol
4  2   20   M   1.67 Lectura
5  3   15   M   1.60 Danca
6  4   19   H   1.76 Bowling
7  5   23   M   1.70 Musica
8  6   18   H   1.65 Lectura
9 > rbind(df, data.frame(Edat=22, Sexe="M", Alcada=1
10   .70, Hobby="Tennis", stringsAsFactors=FALSE))
11   Edat Sexe Alcada   Hobby
12  1   18   H   1.81 Futbol
13  2   20   M   1.67 Lectura
14  3   15   M   1.60 Danca
15  4   19   H   1.76 Bowling
16  5   23   M   1.70 Musica
17  6   18   H   1.65 Lectura
18  7   22   M   1.70 Tennis
19 >

```

Fixeu-vos que, afegint data frame, cal tenir cura de posar l'opció `stringsAsFactors` a `FALSE` si algun dels camps del data frame és una cadena de caràcters. Ara veurem com afegir files mitjançant un list.

```

1 > df
2   Edat Sexe Alcada   Hobby
3  1   18   H   1.81 Futbol
4  2   20   M   1.67 Lectura
5  3   15   M   1.60 Danca
6  4   19   H   1.76 Bowling
7  5   23   M   1.70 Musica

```



```

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

```

Adicionalment, podem afegir més d'una fila, és a dir, combinar data frames:

```

1 > df
2    Edat Sexe Alcada  Hobby
3 1    18    H    1.81  Futbol
4 2    20    M    1.67 Lectura
5 3    15    M    1.60  Danca
6 4    19    H    1.76 Bowling
7 5    23    M    1.70  Musica
8 6    18    H    1.65 Lectura
9 > df1 <- data.frame(Edat=c(17,21,15),Sexe=c("H","M",
    "H"),Alcada=c(1.70,1.71,1.65),Hobby=c("
    Bowling","Danca","Musica"),stringsAsFactors=
    FALSE)
10 > df1
11    Edat Sexe Alcada  Hobby
12 1    17    H    1.70 Bowling
13 2    21    M    1.71  Danca
14 3    15    H    1.65  Musica
15 > rbind(df, df1)
16    Edat Sexe Alcada  Hobby
17 1    18    H    1.81  Futbol
18 2    20    M    1.67 Lectura
19 3    15    M    1.60  Danca
20 4    19    H    1.76 Bowling
21 5    23    M    1.70  Musica
22 6    18    H    1.65 Lectura
23 7    17    H    1.70 Bowling
24 8    21    M    1.71  Danca
25 9    15    H    1.65  Musica
26 >

```

## 2. cbind() per afegir una nova columna (variable).

Per a poder afegir una variable (una columna) al data frame, les dades poden estar guardades directament en un vector, perquè són homogènies

(del mateix tipus). L'única restricció que tenim en aquest cas és que el vector ha de tenir necessàriament el mateix nombre d'elements que files té el 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 > df
5   Edat Sexe Alcada   Hobby   Ciutat
6 1    18   H   1.81 Futbol    Paris
7 2    20   M   1.67 Lectura Barcelona
8 3    15   M   1.60   Danca Barcelona
9 4    19   H   1.76 Bowling    Roma
10 5    23   M   1.70 Musica    Caracas
11 6    18   H   1.65 Lectura Barcelona
12 7    22   M   1.70   Tenis    Paris
13 >

```

3. També podem afegir una columna tal i com fèiem en el list, és a dir fent servir directament l'operador "\$":

```

1 > df
2   Edat Sexe Alcada   Hobby
3 1    18   H   1.81 Futbol
4 2    20   M   1.67 Lectura
5 3    15   M   1.60   Danca
6 4    19   H   1.76 Bowling
7 5    23   M   1.70 Musica
8 6    18   H   1.65 Lectura
9 > df$CodiPostal <- c
    (08027,08003,08014,08034,08034,08006)
10 > df
11   Edat Sexe Alcada   Hobby CodiPostal
12 1    18   H   1.81 Futbol      8027
13 2    20   M   1.67 Lectura      8003
14 3    15   M   1.60   Danca      8014
15 4    19   H   1.76 Bowling      8034
16 5    23   M   1.70 Musica      8034
17 6    18   H   1.65 Lectura      8006
18 >

```

Noteu que, a diferència de la funció `cbind()`, aquesta manera directa d'afegir columnes modifica directament el data frame.

## 4.6 Modificació de data frames

Hi ha alguns problemes on podem necessitar modificar parcialment els data frames. És a dir, podem necessitar modificar els valors d'algunes variables. En aquest cas, podem fer servir una funció de l'R que ens permet fer aquestes

modificacions, `transform()`. A continuació podem veure alguns exemples d'ús d'aquesta funció.

```

1 > df
2   Edat Sexe Alcada   Hobby CodiPostal
3 1   18   H   1.81 Futbol      8027
4 2   20   M   1.67 Lectura      8003
5 3   15   M   1.60  Danca      8014
6 4   19   H   1.76 Bowling      8034
7 5   23   M   1.70 Musica      8034
8 6   18   H   1.65 Lectura      8006
9 > transform(df, Alcada=Alcada*100)
10  Edat Sexe Alcada   Hobby CodiPostal
11 1   18   H  181  Futbol      8027
12 2   20   M  167  Lectura      8003
13 3   15   M  160  Danca      8014
14 4   19   H  176  Bowling      8034
15 5   23   M  170  Musica      8034
16 6   18   H  165  Lectura      8006
17 > transform(df, Alcada=Alcada%%2.54) # passem cm a
    polzades
18   Edat Sexe Alcada   Hobby CodiPostal
19 1   18   H  71.26 Futbol      8027
20 2   20   M  65.75 Lectura      8003
21 3   15   M  62.99  Danca      8014
22 4   19   H  69.29 Bowling      8034
23 5   23   M  66.93 Musica      8034
24 6   18   H  64.96 Lectura      8006
25 > transform(df, Alcada=Alcada%%12, CodiPostal=
    CodiPostal%%100) # passem polzades a peus i reduim
    cp
26   Edat Sexe Alcada   Hobby CodiPostal
27 1   18   H   5.94 Futbol      27
28 2   20   M   5.48 Lectura      3
29 3   15   M   5.25  Danca      14
30 4   19   H   5.77 Bowling      34
31 5   23   M   5.58 Musica      34
32 6   18   H   5.41 Lectura      6
33 >

```

## 4.7 Ordenació de data frames

Els data frames poden ser ordenats per una o més d'una columna. Per exemple, podem necessitar ordenar el data frame `df` segons la variable `Edat`, en ordre creixent. Per a fer això cal fer servir l'operador `order` con segueix:

```

1 > df
2   Edat Sexe Alcada   Hobby CodiPostal
3 1   18   H   1.81 Futbol      8027
4 2   20   M   1.67 Lectura      8003

```

```

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

```

També podem ordenar el data frame en ordre decreixent posant el flag `decreasing` a `TRUE` en la funció `order` (per defecte és `FALSE`):

```

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

```

I ara, creixent per `Edat` i per `Alcada`

```

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

```

```

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

```

## 4.8 Fusió de data frames

`merge(x,y)`: A l'R és possible construir data frames a partir de la fusió de dos data frames `x` i `y`. Els data frames es combinen segons els valors d'una variable (columna) comú. El data frame resultant de fusionar dos data frames tindrà les columnes dels dos data frames -sense repetir les que siguin comuns- i només les files d'ambdós data frames on els valors de la columna en comú es correspon en els dos data frames:

```

1 > D1
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 > D2
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(D1,D2)
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 volem el resultat del `merge` amb totes les files malgrat que hi hagi valors de la columna compartida que no siguin comuns. Fixeu-vos que a les files on hi ha valors no compartits en la columna `Var2` es coloca `<NA>` en els espais on no existeix valor a col·locar:

```

1  > D1
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 > D2
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(D1,D2,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 >

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  > D1
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 > D2
13 >
14    Dif2 Var4
15 1    B1   D1
16 2    B2   D2
17 3    B3   D1
18 4    B4   D3
19 5    B5   D5
20 >
21 > merge(D1,D2,by.x = "Var2",by.y = "Dif2")
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 >
```