

Programación Estadística: Manejo de Datos

Jocelyn Simmonds (jsimmond@dcc.uchile.cl)

Departamento de Ciencias de la Computación

Se puede usar la función `which()` para encontrar las posiciones de los elementos de un vector que cumplen una condición:

```
1 > x <- 1:10
2 > which(x %% 2 == 0)
3 [1] 2 4 6 8 10
4 >
5 > x2 <- runif(10, 2, 4.5) # genera numeros al azar (dist. uniforme)
6 > x2
7 [1] 3.966691 4.420197 3.243075 3.012766 3.617560 3.299832 2.524505 2.098100
8 [9] 2.113189 4.041362
9 > which(x2 >= mean(x2))
10 [1] 1 2 3 5 6 10
```

Se usan las comillas para crear secuencias de caracteres:

```
1 > x <- "hola"
2 > x
3 [1] "hola"
4 > y <- 'chao'
5 > y
6 [1] "chao"
7 > z <- "hola y chao"
8 > z
9 [1] "hola y chao"
10 > c(x, "y", y)
11 [1] "hola" "y"      "chao"
12 > paste(x, "y", y)
13 [1] "hola y chao"
```

La función `paste()` se puede usar para generar los rótulos de un gráfico:

```
1 > labels <- paste("Libro", 1:10)
2 > labels
3 [1] "Libro 1" "Libro 2" "Libro 3" "Libro 4" "Libro 5" "Libro 6"
4 [7] "Libro 7" "Libro 8" "Libro 9" "Libro 10"
5 > labels2 <- paste("x", 1:10, 10:4)
6 > labels2
7 [1] "x 1 10" "x 2 9" "x 3 8" "x 4 7" "x 5 6" "x 6 5" "x 7 4" "x 8 10"
8 [9] "x 9 9" "x 10 8"
9 > labels3 <- paste("Lab", 1:10, sep="")
10 > labels3
11 [1] "Lab1" "Lab2" "Lab3" "Lab4" "Lab5" "Lab6" "Lab7" "Lab8" "Lab9"
12 [10] "Lab10"
```

Si alguno de los argumentos es más corto que el resto, se “repiten” los elementos según sea necesario.

Creen los siguientes vectores:

- 1 "x 1" "y 2" "x 3" "y 4" "x 5"
- 2 "x 1 10" "y 2 9" "x 3 8" "y 4 7" "x 5 10"

Objetos

Tipos de objetos

- vector: secuencia de elementos, todos del mismo tipo
- factor: vector especial para manejar datos categóricos
- arreglo: colección de elementos de datos, organizados en varias dimensiones (p.ej., cubo de $3 \times 3 \times 3$)
- matriz: arreglo de dos dimensiones
- lista: secuencia de elementos, pueden ser de cualquier tipo (incluyendo listas)
- data frame: lista de vectores del mismo largo

Las funciones también son objetos en R, pero las veremos en la siguiente clase.

Propiedades de objetos

Todo objeto x tiene al menos dos propiedades:

- 1 `mode(x)`: indica el “tipo” del objeto. Ejemplos: `numeric`, `logical`, `character`, etc.
- 2 `length(x)`: indica el “tamaño” del objeto (no siempre es útil)

```
1 > x <- 1:10
2 > mode(x)
3 [1] "numeric"
4 > length(x)
5 [1] 10
6 >
7 > x <- "hola"
8 > mode(x)
9 [1] "character"
10 > length(x)
11 [1] 1
```

Cambio de mode()

Existen varias funciones `as.xxxxx(objeto)`, que permiten cambiar el `mode()` de un objeto.

```
1 > z <- 0:9
2 > digitos <- as.character(z)
3 > digitos
4 [1] "0" "1" "2" "3" "4" "5" "6" "7" "8" "9"
5 >
6 > d <- as.integer(digitos)
7 > d
8 [1] 0 1 2 3 4 5 6 7 8 9
9 > d == z
10 [1] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
```

También existen varias funciones `is.xxxxx(objeto)`, que retornan `TRUE` si el objeto es del tipo en el nombre de la función.

Pruébenlo

- 1 Creen un vector de números decimales, vean que pasa cuando aplican `as.integer()`
- 2 ¿Qué pasa si le aplican `as.integer()` a un vector de texto?
- 3 ...y a un vector de valores lógicos?

Vectores

Cambiando el tamaño de un vector

Podemos asignar valores en posiciones mayores que el largo actual:

```
1 > x <- integer(15)
2 > x
3 [1] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
4 >
5 > e <- integer()      # vector de largo 0
6 > e[3] <- 17
7 > e
8 [1] NA NA 17          # ojo que rellena con NA en vez de 0
9 > e[10] <- 3
10 > e
11 [1] NA NA 17 NA NA NA NA NA NA 3
```

Cambiando elementos del vector

Podemos cambiar elementos usando el complemento:

```
1 > e
2 [1] NA NA 17 NA NA NA NA NA NA 3
3 > e[-2] <- 4 # asignar 4 a todo menos la posición 2
4 > e
5 [1] 4 NA 4 4 4 4 4 4 4 4
6 > e[2] <- 1.3
7 > e
8 [1] 4.0 1.3 4.0 4.0 4.0 4.0 4.0 4.0 4.0 4.0
9 > mode(e)
10 [1] "numeric"
11 > e[c(1, 5, 8)] <- "abc"
12 > e
13 [1] "abc" "1.3" "4" "4" "abc" "4" "4" "abc" "4" "4"
14 >
```

Noten que el vector pasa a ser de tipo `numeric` por la asignación en la línea 6, y `character` por la asignación en la línea 11.

Cambiando el tamaño de un vector

También podemos cambiar el valor de `length()`:

```
1 > f
2 [1] 4 4 NA NA NA
3 >
4 > length(f) <- 3
5 > f
6 [1] 4 4 NA
7 >
8 > length(f) <- 13
9 > f
10 [1] 4 4 NA NA NA NA NA NA NA NA NA NA NA
```

En general, no recomiendo hacer esto ...

Concatenación de elementos

Podemos agregar nuevos elementos a la cola del vector usando `c()`:

```
1 > x <- NULL
2 > x <- c(x, 2)
3 > x
4 [1] 2
5 > x <- c(x, 4)
6 > x
7 [1] 2 4
8 > x <- c(x, 6)
9 > x
10 [1] 2 4 6
11 > x <- c(x, 8)
12 > x
13 [1] 2 4 6 8
14 >
```

Creen el vector x de largo 15:

- 1 asignen el valor 1.4 a las posiciones 2, 7 y 9
- 2 asignen el valor 2.8 a las posiciones 4, 8 y 12
- 3 ... y 1 al resto.

Factores

Factores

Podemos crear un factor de un vector de datos categóricos:

```
1 > x <- c("rojo", "azul", "rojo", "rojo", "azul")
2 > fact_x <- factor(x)
3 > fact_x
4 [1] rojo azul rojo rojo azul
5 Levels: azul rojo          # siempre se muestran en orden lexicografico
6 >
7 > levels(fact_x)
8 [1] "azul" "rojo"
9 >
10 > y <- c("rojo1", "azul", "rojo2", "rojo11", "azul")
11 > fact_y <-factor(y)
12 > fact_y
13 [1] rojo1 azul rojo2 rojo11 azul
14 Levels: azul rojo1 rojo11 rojo2
```

Los Levels de un factor son las categorías que R identificó en los datos.

Factores: uso

Los factores se usan para calcular funciones como el promedio o suma en forma agregada por categoría:

```
1 > fact_x
2 [1] rojo azul rojo rojo azul
3 Levels: azul rojo
4 > datos <- 1:5
5 > tapply(datos, fact_x, sum)
6 azul rojo
7     7     8
8 > tapply(datos, fact_x, mean)
9     azul     rojo
10 3.500000 2.666667
11 > tapply(datos, fact_x, max)
12 azul rojo
13     5     4
```

Ojo: datos y el vector usado para crear fact_x debe ser del mismo largo.

Combinando factores

La función `tapply()` recibe como segundo parámetro una lista de factores:

```
1 > fact_x
2 [1] rojo azul rojo rojo azul
3 Levels: azul rojo
4 > datos <- 1:5
5 >
6 > genero <- c("H", "H", "F", "F", "F")
7 > fact_g <- factor(genero)
8 > tapply(datos, list(fact_x, fact_g), sum)
9      F H
10 azul 5 2
11 rojo 7 1
```

Pruébenlo

Creemos de nuevo el vector de edades de los compañeros del curso.
Colectemos además los siguientes datos:

- 1 preferencias de películas: drama, comedia, acción o misterio?
- 2 sabor favorito de helado: frutilla, vainilla o chocolate?

Hagan el cruce de datos usando `tapply()`

Arreglos y matrices

Creando un arreglo

```
1 > dim(z) <- c(3,5,100)
2 > z
3 ...
4 , , 99
5
6      [,1] [,2] [,3] [,4] [,5]
7 [1,]    0    0    0    0    0
8 [2,]    0    0    0    0    0
9 [3,]    0    0    0    0    0
10
11 , , 100
12
13      [,1] [,2] [,3] [,4] [,5]
14 [1,]    0    0    0    0    0
15 [2,]    0    0    0    0    0
16 [3,]    0    0    0    0    0
17 > dim(z) <- c(3,6,100)
18 Error in dim(z) <- c(3, 6, 100) :
19   dims [product 1800] do not match the length of object [1500]
```


Creando un arreglo

Puede redefinir las dimensiones si es consistente con número de elementos:

```
1 > dim(z) <- c(30,50)
2 > z
3      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10] [,11] [,12] [,13]
4 [1,]    0    0    0    0    0    0    0    0    0    0    0    0    0
5 [2,]    0    0    0    0    0    0    0    0    0    0    0    0    0
6 ...
7 [30,]   0    0    0    0    0    0    0    0    0    0    0    0    0
8 ...
9 ...
10     [,38] [,39] [,40] [,41] [,42] [,43] [,44] [,45] [,46] [,47] [,48] [,49]
11 [1,]     0    0    0    0    0    0    0    0    0    0    0    0
12 [2,]     0    0    0    0    0    0    0    0    0    0    0    0
13 ...
14 [30,]    0    0    0    0    0    0    0    0    0    0    0    0
15     [,50]
16 [1,]     0
17 [2,]     0
18 ...
19 [30,]    0
```

Orden de elementos

Los elementos de un arreglo están ordenados por columnas. En otras palabras, la matriz

$$M = \begin{bmatrix} 11 & 12 & 13 \\ 21 & 22 & 23 \end{bmatrix}$$

se crea de la siguiente forma:

```
1 > M <- c(11, 21, 12, 22, 13, 23)
2 > dim(M) <- c(2, 3)
3 > M
4      [,1] [,2] [,3]
5 [1,]  11  12  13
6 [2,]  21  22  23
7 >
```

Acceder a elementos por posición

Usaremos `[]` para acceder a los elementos de un arreglo por posición:

```
1 > z <- numeric(1500)
2 > dim(z) <- c(3,5,100)
3 > z[1,1,99] <- 1
4 > z[, ,99] # se muestra la matriz completa del nivel 99
5      [,1] [,2] [,3] [,4] [,5]
6 [1,]    1    0    0    0    0
7 [2,]    0    0    0    0    0
8 [3,]    0    0    0    0    0
```

Si omitimos una posición, R incluye la columna/fila completa.

Acceder a elementos por posición

```
1 > z[1,1,99] <- 1
2 > z[, ,99]
3     [,1] [,2] [,3] [,4] [,5]
4 [1,]    1    0    0    0    0
5 [2,]    0    0    0    0    0
6 [3,]    0    0    0    0    0
7 > z[1,3,99] <- 3
8 > z[, ,99]
9     [,1] [,2] [,3] [,4] [,5]
10 [1,]    1    0    3    0    0
11 [2,]    0    0    0    0    0
12 [3,]    0    0    0    0    0
13 > z[1, ,98:100]
14     [,1] [,2] [,3]
15 [1,]    0    1    0
16 [2,]    0    0    0
17 [3,]    0    3    0
18 [4,]    0    0    0
19 [5,]    0    0    0
20 >
```

Acceder a elementos por posición

Una vez que hayamos indicado las dimensiones de un arreglo, podemos usar `dim()` para ver las dimensiones actuales del arreglo:

```
1 > dim(z)
2 [1] 3 5 100
```

De todas formas, podemos acceder a los elementos del arreglo directamente por su posición en el vector de datos:

```
1 > z[1,1,99] <- 1
2 > z[1]
3 [1] 0
4 > z[1470]      # 3 x 5 x 98
5 [1] 0
6 > z[1471]
7 [1] 1
```

Otra forma de crear un arreglo

También podemos crear un arreglo en un solo paso, usando la función `array(vector de datos, dim = vector de dimensiones)`:

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

Ojo: si el vector de datos es más corto que el número de elementos especificados por las dimensiones, el vector de datos se “reusa” hasta completar el arreglo.

Arreglos de posiciones

Podemos crear un “arreglo de posiciones” para modificar un arreglo:

```
1 > x <- array(1:20, dim=c(4,5))
2 > i <- array(c(1, 2, 3, 4, 1, 3), dim=c(3, 2))
3 > i
4      [,1] [,2]
5 [1,]    1    4
6 [2,]    2    1
7 [3,]    3    3
8 > x[i] <- 0
9 > x
10     [,1] [,2] [,3] [,4] [,5]
11 [1,]    1    5    9    0   17
12 [2,]    0    6   10   14   18
13 [3,]    3    7    0   15   19
14 [4,]    4    8   12   16   20
```

Ojo: posiciones negativas generan un error; posiciones con algún 0 son ignorados.

Matrices

Matrices son arreglos de 2 dimensiones, pueden ser incluidas en expresiones aritméticas:

```
1 > C <- array(c(4, 3, 3, 2), dim=c(2,2))
2 > C
3      [,1] [,2]
4 [1,]    4    3
5 [2,]    3    2
6 >
7 > C + 2
8      [,1] [,2]
9 [1,]    6    5
10 [2,]    5    4
11 >
12 > C * C      # ojo! * no es multiplicacion de matrices!
13      [,1] [,2]
14 [1,]   16    9
15 [2,]    9    4
```

Operaciones sobre matrices

Sean A , D matrices, b , x vectores y n una constante:

- $t(A)$: retorna A traspuesta
- $nrow(A)$: retorna el número de filas de A
- $ncol(A)$: retorna el número de columnas de A
- $A \%*\% B$: retorna el producto de las matrices A , B
- $det(A)$: retorna el determinante de A
- $diag(A)$: retorna un vector con la diagonal principal de A
- $diag(x)$: retorna una matriz diagonal, donde los elementos de la diagonal son los valores de x
- $diag(n)$: retorna una matriz identidad de $n \times n$
- $solve(A, b)$: resuelve el sistema de ecuaciones lineales $A \cdot x = b$
- $solve(A)$: retorna la matriz inversa de A

Resuelvan el siguiente sistema de ecuaciones lineales:

$$2x + y - 2z = 10$$

$$3x + 2y + 2z = 1$$

$$5x + 4y + 3z = 4$$

Listas y Data frames

Listas: elementos de cualquier tipo

ex es una lista de tres elementos: uno de tipo character, uno de tipo integer y uno de tipo vector

```
1 > ex <- list(nombre="Juan", num.hijos=2, edades.hijos=c(4, 7))
2 > ex
3 $nombre
4 [1] "Juan"
5
6 $num.hijos
7 [1] 2
8
9 $edades.hijos
10 [1] 4 7
```

Listas: elementos de cualquier tipo

```
1 > ex <- list(nombre="Juan", num.hijos=2, edades.hijos=c(4, 7))
2 > ex[3]
3 $edades.hijos
4 [1] 4 7
5
6 > ex[["num.hijos"]]
7 [1] 2
8 >
9 > ex[[3]]
10 [1] 4 7
11 >
12 > ex[[3]][1]
13 [1] 4
14 >
15 > ex$edades.hijos
16 [1] 4 7
17 >
18 > ex$edades.hijos[2]
19 [1] 7
```

Listas: elementos de cualquier tipo

R no revisa que los nombres de los elementos de una lista sean únicos, lo que puede causar errores a futuro:

```
1 > ex2 <- list(x=1, y=2, x=3)
2 > ex2
3 $x
4 [1] 1
5
6 $y
7 [1] 2
8
9 $x
10 [1] 3
11
12 > ex2[["x"]] # toma el primer elemento con nombre = x
13 [1] 1
14 > ex2[[1]]
15 [1] 1
16 > ex2[[3]]
17 [1] 3
```

Editando una listas

Pueden usar `c()` para unir listas:

```
1 > ex3 <- c(ex, ex2)
2 > ex3
3 $nombre
4 [1] "Juan"
5
6 $num.hijos
7 [1] 2
8
9 $edades.hijos
10 [1] 4 7
11
12 $x
13 [1] 1
14
15 $y
16 [1] 2
17
18 $x
19 [1] 3
```

Editando una listas

Tambien pueden usar `length()` para acortar una lista:

```
1 > length(ex3)
2 [1] 6
3 > length(ex3) <- 3
4 > ex3
5 $nombre
6 [1] "Juan"
7
8 $num.hijos
9 [1] 2
10
11 $edades.hijos
12 [1] 4 7
```

Data frame

Un data frame es una lista donde todas las columnas son del mismo largo:

```
1 > mtcars
2           mpg cyl  disp  hp drat   wt  qsec vs am gear carb
3 Mazda RX4      21.0   6  160.0 110 3.90 2.620 16.46 0  1   4   4
4 Mazda RX4 Wag  21.0   6  160.0 110 3.90 2.875 17.02 0  1   4   4
5 Datsun 710     22.8   4  108.0  93 3.85 2.320 18.61 1  1   4   1
6 Hornet 4 Drive 21.4   6  258.0 110 3.08 3.215 19.44 1  0   3   1
7 Hornet Sportabout 18.7   8  360.0 175 3.15 3.440 17.02 0  0   3   2
8 Valiant        18.1   6  225.0 105 2.76 3.460 20.22 1  0   3   1
9 ...
```

Un data frame incluye nombres de filas y nombres de columnas.

R trae varios data frames de prueba, pueden listarlos usando la función `data()`.

Creando un data frame

Pueden crear un data frame de las siguientes formas:

- de forma manual, usando la función `data.frame()`,
- convirtiendo una lista de listas, usando la función `as.data.frame()`,
- ...o cargando los datos de un archivo externo, usando la función `read.table()`

Leer datos de un archivo de texto

Para usar `read.table()`, los datos deben estar en el siguiente formato:

- la primera fila del archivo tiene los nombres de las n columnas
- el resto de las filas contienen los datos de los registros ($n + 1$ datos porque incluye un identificador de fila)

Ejemplo: `houses.data`

```
1 Price Floor Area Rooms Age Cent.heat
2 01 52.00 111.0 830 5 6.2 no
3 02 54.75 128.0 710 5 7.5 no
4 03 57.50 101.0 1000 5 4.2 no
5 04 57.50 131.0 690 6 8.8 no
6 05 59.75 93.0 900 5 1.9 yes
```

Ahora podemos cargar los datos:

```
1 > # archivo debe estar dentro del workspace
2 > prices <- read.table("houses.data")
```

Acceso a los componentes de un data frame

Usaremos el signo \$ para indicar la componente que nos interesa:

```
1 > prices <- read.table("houses.data")
2 > prices
3   Price Floor Area Rooms Age Cent.heat
4 01 52.00   111  830     5 6.2         no
5 02 54.75   128  710     5 7.5         no
6 03 57.50   101 1000     5 4.2         no
7 04 57.50   131  690     6 8.8         no
8 05 59.75    93  900     5 1.9         yes
9 > prices$Floor
10 [1] 111 128 101 131  93
11 > prices[[1]]
12 [1] 52.00 54.75 57.50 57.50 59.75
13 > prices["03", "Rooms"]
14 [1] 5
15 > prices["03", 1:5]
16   Price Floor Area Rooms Age
17 03  57.5   101 1000     5 4.2
```

Editando un data frame

R provee una interfaz para editar los elementos de un data frame:

```
1 > prices_new <-edit(prices)
```

```
> prices
  Price Floor Area Rooms Age Cent.heat
01 52.00  111  830    5 6.2         no
02 54.75  128  710    5 7.5         no
03 57.50  101 1000    5 4.2         no
04 57.50  131  690    6 8.8         no
05 59.75   93  900    5 1.9         no
> prices$Floor
[1] 111 128 101 131  93
> prices[[1]]
[1] 52.00 54.75 57.50 57.50 59.75
> prices["03", Rooms]
Error in `[.data.frame'](prices,
> prices["03", "Rooms"]
[1] 5
> prices["03", 1:5]
  Price Floor Area Rooms Age
03 57.5  101 1000    5 4.2
> prices_new <-edit(prices)
```

	row.names	Price	Floor	Area	Rooms	Age	Cent.heat
1	01	52	111	830	5	6.2	no
2	02	54.75	128	710	5	7.5	no
3	03	57.5	101	1000	5	4.2	no
4	04	57.5	131	690	6	8.8	no
5	05	59.75	93	900	5	1.9	yes
6	06	45	76	680	4	1.5	yes
7							
8							
9							
10							
11							
12							

Editando un data frame

Pueden editar el data frame original en forma directa:

```
1 > fix(prices)
```

O pueden crear un nuevo data frame:

```
1 > nuevo <- edit(data.frame())
```

Pruébenlo

Usen los datos de películas y sabores de helados para crear su propio archivo de datos. Carguen estos datos en un data frame usando `read.table()`

Creando un data frame en forma manual

Pueden crear un data frame a partir de vectores existentes:

```
1 > nombres <- c("Ana Maria", "Doris", "Paula", "Pedro", "Juan")
2 > apellidos <- c("Jerez", "Said", "Soto", "Perez", "Soto")
3 > edades <- c(23, 65, 44, 45, 33)
4 > personas <- as.data.frame(list(nombres, apellidos, edades))
5 > personas
6   c..Ana.Maria....Doris....Paula....Pedro....Juan..
7 1                               Ana Maria
8 2                               Doris
9 3                               Paula
10 4                               Pedro
11 5                               Juan
12  c...Jerez....Said....Soto....Perez....Soto.. c.23..65..44..45..33.
13 1                               Jerez                23
14 2                               Said                 65
15 3                               Soto                 44
16 4                               Perez                45
17 5                               Soto                 33
```

Falta crear los nombres de las componentes.

Creando un data frame en forma manual

Pueden usar `names()` para definir los nombres de las componentes de un data frame:

```
1 > names(personas) <- c("Nombre", "Apellido", "Edad")
2 > personas
3     Nombre Apellido Edad
4 1 Ana Maria   Jerez   23
5 2   Doris    Said    65
6 3   Paula    Soto    44
7 4   Pedro   Perez    45
8 5    Juan    Soto    33
9 > personas$Apellido
10 [1] Jerez Said  Soto  Perez Soto
11 Levels: Jerez Perez Said Soto
12 > personas$Edad
13 [1] 23 65 44 45 33
```
