

# Notas sobre R:

## Un entorno de programación para Análisis de Datos y Gráficos

Bill Venables & Dave Smith

*Department of Statistics  
The University of Adelaide*

Robert Gentleman & Ross Ihaka

*Department of Statistics  
University of Auckland*

Martin Mächler

*Statistics Seminar  
ETH Zurich*

Andrés González y Silvia González

*Instituto de Estadística de Andalucía*

© W. Venables, 1990, 1992.

© R. Gentleman & R. Ihaka, 1997.

© M. Mächler, 1997–1998.

© Traducción español, A. González y S. González, 2000.

# Prólogo

Estas notas sobre R están escritas a partir de un conjunto de notas que describían los entornos S y S-PLUS escritas por Bill Venables y Dave Smith. Hemos realizado un pequeño número de cambios para reflejar las diferencias entre R y S.

R es un proyecto vivo y sus capacidades no coinciden totalmente con las de S. En estas notas hemos adoptado la convención de que cualquier característica que se vaya a implementar se especifica como tal en el comienzo de la sección en que la característica es descrita. Los usuarios pueden contribuir al proyecto implementando cualquiera de ellas.

Deseamos dar las gracias más efusivas a Bill Venables por permitir la distribución de esta versión modificada de las notas y por ser un defensor de R desde su inicio.

Cualquier comentario o corrección serán siempre bienvenidos. Dirija cualquier correspondencia a

`R@stat.auckland.ac.nz`.

## Sugerencias al lector

La primera relación con R debería ser la sesión inicial del Apéndice A. Está escrita para que se pueda conseguir cierta familiaridad con el estilo de las sesiones de R y para comprobar que coincide con la versión actual.

Muchos usuarios eligen R fundamentalmente por sus capacidades gráficas. Si ese es su caso, debería leer antes o después la sección 11 sobre capacidades gráficas y para ello no es necesario esperar a haber asimilado totalmente las secciones precedentes.

Robert Gentleman and Ross Ihaka,  
University of Auckland,  
Abril de 1997.

# Índice General

## Prólogo

<b>1</b>	<b>Introducción y Preliminares</b>	<b>2</b>
1.1	El entorno R . . . . .	2
1.2	Programas relacionados y Documentación . . . . .	2
1.3	R y el sistema MICROSOFT WINDOWS . . . . .	2
1.4	Utilización interactiva de R . . . . .	3
1.5	Una sesión inicial . . . . .	3
1.6	Ayuda sobre funciones y capacidades . . . . .	4
1.7	Órdenes de R. Mayúsculas y minúsculas . . . . .	4
1.8	Recuperación y corrección de órdenes previas . . . . .	4
1.9	Ejecución de órdenes desde un archivo y redirección de la salida . . . . .	4
1.10	Almacenamiento y eliminación de objetos . . . . .	5
<b>2</b>	<b>Cálculos sencillos. Números y vectores</b>	<b>6</b>
2.1	Vectores (numéricos). Asignación . . . . .	6
2.2	Aritmética vectorial . . . . .	6
2.3	Generación de sucesiones . . . . .	7
2.4	Vectores lógicos . . . . .	8
2.5	Valores faltantes . . . . .	8
2.6	Vectores de caracteres . . . . .	9
2.7	Vectores de índices. Selección y modificación de subvectores . . . . .	9
<b>3</b>	<b>Objetos: Modos y atributos</b>	<b>11</b>
3.1	Atributos intrínsecos: <i>modo</i> y <i>longitud</i> . . . . .	11
3.2	Modificación de la longitud de un objeto . . . . .	11
3.3	<code>attributes()</code> y <code>attr()</code> . . . . .	12
3.4	Clases de objetos . . . . .	12
<b>4</b>	<b>Factores Nominales y Ordinales</b>	<b>14</b>
4.1	Un ejemplo específico . . . . .	14
4.2	La función <code>tapply()</code> . Variables desastradas (ragged arrays) . . . . .	14
<b>5</b>	<b>Variables indexadas. Matrices</b>	<b>16</b>
5.1	Variables indexadas (Arrays) . . . . .	16
5.2	Elementos de una variable indexada . . . . .	16
5.3	Uso de variables indexadas como índices . . . . .	17
5.4	La función <code>array()</code> . . . . .	18
5.4.1	Operaciones con variables indexadas y vectores. Reciclado. . . . .	18
5.5	Producto exterior de dos variables indexadas . . . . .	19
5.5.1	Ejemplo: Distribución del determinante de una matriz $2 \times 2$ de dígitos . . . . .	19
5.6	Traspuesta generalizada de una variable indexada . . . . .	20
5.7	Operaciones con matrices: Producto matricial. Inversa de una matriz. Resolución de sistemas lineales . . . . .	20
5.8	Submatrices. Funciones <code>cbind()</code> y <code>rbind()</code> . . . . .	21
5.9	La función de concatenación, <code>c()</code> , con variables indexadas . . . . .	21
5.10	Tablas de frecuencias a partir de factores. La función <code>table()</code> . . . . .	21
<b>6</b>	<b>Listas y hojas de datos. Utilización</b>	<b>23</b>
6.1	Listas . . . . .	23
6.2	Construcción y modificación de listas . . . . .	24
6.2.1	Concatenación de listas . . . . .	24
6.3	Algunas funciones que devuelven una lista . . . . .	24
6.3.1	Autovalores y autovectores . . . . .	24

6.3.2	Descomposición en valores singulares. Determinantes . . . . .	25
6.3.3	Ajuste por mínimos cuadrados. Descomposición <i>QR</i> . . . . .	25
6.4	Hojas de datos (Data frames) . . . . .	26
6.4.1	Construcción de hojas de datos . . . . .	26
6.4.2	Funciones <code>attach()</code> y <code>detach()</code> . . . . .	26
6.4.3	Trabajo con hojas de datos . . . . .	27
<b>7</b>	<b>Lectura de datos de un archivo</b>	<b>28</b>
7.1	La función <code>read.table()</code> . . . . .	28
7.2	La función <code>scan()</code> . . . . .	29
<b>8</b>	<b>Ciclos. Ejecución condicional</b>	<b>30</b>
8.1	Expresiones agrupadas . . . . .	30
8.2	Órdenes de control . . . . .	30
8.2.1	Ejecución condicional: la función <code>if</code> . . . . .	30
8.2.2	Ciclos: Funciones <code>for</code> , <code>repeat</code> y <code>while</code> . . . . .	30
<b>9</b>	<b>Escritura de nuevas funciones</b>	<b>32</b>
9.1	Ejemplos elementales . . . . .	32
9.2	Cómo definir un operador binario . . . . .	33
9.3	Argumentos con nombre. Valores predeterminados. Argumento “...” . . . . .	33
9.4	Las asignaciones dentro de una función son locales. Marco. . . . .	34
9.5	Ejemplos más complejos . . . . .	34
9.5.1	Factores de eficiencia en diseño en bloques . . . . .	34
9.5.2	Cómo eliminar los nombres al imprimir una variable indexada . . . . .	35
9.5.3	Integración numérica recursiva . . . . .	35
9.6	Ámbito . . . . .	37
9.7	Adaptación del entorno . . . . .	38
9.8	Clases. Funciones genéricas. Orientación a objetos . . . . .	39
<b>10</b>	<b>Modelos estadísticos en R</b>	<b>41</b>
10.1	Definición de modelos estadísticos. Fórmulas . . . . .	41
10.2	Modelos de regresión . . . . .	43
10.3	Funciones genéricas de extracción de información . . . . .	43
10.4	Análisis de varianza. Comparación de modelos . . . . .	43
10.4.1	Tablas ANOVA . . . . .	44
10.5	Actualización de modelos ajustados. Uso de “.” . . . .	45
10.6	Modelos lineales generalizados. Familias . . . . .	45
10.6.1	Familias . . . . .	46
10.6.2	La función <code>glm</code> . . . . .	47
10.7	Modelos de Mínimos cuadrados no lineales y de Máxima verosimilitud . . . . .	49
10.7.1	Mínimos cuadrados . . . . .	49
10.7.2	Máxima verosimilitud . . . . .	50
10.8	Algunos modelos no-estándar . . . . .	51
<b>11</b>	<b>Procedimientos gráficos</b>	<b>52</b>
11.1	Funciones gráficas de nivel alto . . . . .	52
11.1.1	La función <code>plot</code> . . . . .	52
11.1.2	Representación de datos multivariantes . . . . .	53
11.1.3	Otras representaciones gráficas . . . . .	53
11.1.4	Argumentos de las funciones gráficas de nivel alto . . . . .	54
11.2	Funciones gráficas de nivel bajo . . . . .	55
11.2.1	Anotaciones matemáticas . . . . .	56
11.2.2	Fuentes vectoriales Hershey . . . . .	56
11.3	Funciones gráficas interactivas . . . . .	57
11.4	Uso de parámetros gráficos . . . . .	57

11.4.1	Cambios permanentes. La función <code>par()</code> . . . . .	58
11.4.2	Cambios temporales. Argumentos de las funciones gráficas . . . . .	58
11.5	Parámetros gráficos habituales . . . . .	59
11.5.1	Elementos gráficos . . . . .	59
11.5.2	Ejes y marcas de división . . . . .	60
11.5.3	Márgenes de las figuras . . . . .	61
11.5.4	Figuras múltiples . . . . .	62
11.6	Dispositivos gráficos . . . . .	63
11.6.1	Inclusión de gráficos <code>PostScript</code> en documentos . . . . .	64
11.6.2	Dispositivos gráficos múltiples . . . . .	64
<b>A</b>	<b>Primera sesión con R</b>	<b>66</b>
<b>B</b>	<b>El editor de órdenes de R</b>	<b>69</b>
B.1	Preliminares . . . . .	69
B.2	Edición de acciones . . . . .	69
B.3	Resumen del editor de líneas de órdenes . . . . .	69

# 1 Introducción y Preliminares

## 1.1 El entorno R

R es un conjunto integrado de programas para manipulación de datos, cálculo y gráficos. Entre otras características dispone de:

- almacenamiento y manipulación efectiva de datos,
- operadores para cálculo sobre variables indexadas (arrays), en particular matrices,
- una amplia, coherente e integrada colección de herramientas para análisis de datos,
- posibilidades gráficas para análisis de datos, que funcionan directamente sobre pantalla o impresora, y
- un lenguaje de programación bien desarrollado, simple y efectivo, que incluye condicionales, ciclos, funciones recursivas y posibilidad de entradas y salidas. (Debe destacarse que muchas de las funciones suministradas con el sistema están escritas en el lenguaje S).

El término “entorno” lo caracteriza como un sistema completamente diseñado y coherente, antes que como una agregación incremental de herramientas muy específicas e inflexibles, como es frecuentemente el caso con otros programas de análisis de datos.

R es en gran parte un vehículo para el desarrollo de nuevos métodos de análisis interactivo de datos. Como tal es muy dinámico y las diferentes versiones no siempre son totalmente compatibles con las anteriores. Algunos usuarios prefieren los cambios debido a los nuevos métodos y tecnología que los acompañan, a otros sin embargo les molesta ya que algún código anterior deja de funcionar. Aunque R puede entenderse como un lenguaje de programación, los programas escritos en R deben considerarse esencialmente efímeros.

## 1.2 Programas relacionados y Documentación

R puede entenderse como una nueva implementación del lenguaje S desarrollado en AT&T por Rick Becker, John Chambers y Allan Wilks. Muchos de los libros y manuales sobre S son útiles para R.

La referencia básica es Richard A. Becker, John M. Chambers and Allan R. Wilks (1988) de Richard A. Becker, John M. Chambers and Allan R. Wilks. Las características de la versión de agosto de 1991 de S están recogidas en WhiteBook editado por John M. Chambers y Trevor J. Hastie. Además, los manuales de S-PLUS, la versión comercial de S, pueden ser útiles.

## 1.3 R y el sistema MICROSOFT WINDOWS

La forma más conveniente de usar R es en una estación de trabajo con un sistema de ventanas. Estas notas están escritas pensando en usuarios de estas características. En particular nos referiremos ocasionalmente a la utilización de R en un sistema X-windows, aunque normalmente se pueden aplicar a cualquier implementación del entorno R.

Muchos usuarios encontrarán necesario interactuar directamente con el sistema operativo de su ordenador de vez en cuando. En estas notas se trata fundamentalmente de la interacción con el sistema operativo UNIX. Si utiliza R bajo MACINTOSH o MICROSOFT WINDOWS necesitará realizar algunos pequeños cambios. En ambas implementaciones hemos tratado de atenarnos a

El ajuste del ordenador para obtener el máximo rendimiento de las cualidades parametrizables de R es una tarea interesante aunque tediosa y no se considerará en estas notas. Si tiene dificultades busque a un experto cercano a usted.

## 1.4 Utilización interactiva de R

Al utilizar R éste presenta un símbolo cuando espera la entrada de órdenes. El símbolo predeterminado es “>”, que en UNIX puede coincidir con el símbolo del sistema, por lo que puede parecer que no sucede nada. Por tanto, es posible modificar este símbolo en R. En estas notas supondremos que el símbolo de UNIX es “\$”.

Para utilizar R bajo UNIX por primera vez, el procedimiento recomendado es el siguiente:

1. Cree un subdirectorio, por ejemplo **trabajo**, en el que crear los archivos de datos que desee analizar mediante R. Éste será el directorio de trabajo cada vez que utilice R para este problema concreto.

```
$ mkdir trabajo
```

```
$ cd trabajo
```

2. Inicie R con la orden

```
$ R
```

3. Ahora puede escribir órdenes para R (como se hace más adelante).

4. Para salir de R la orden es

```
> q()
```

```
$
```

R preguntará si desea salvar los datos de esta sesión de trabajo. Puede responder **yes** (Si), **no** (No) o **cancel** (cancelar) pulsando respectivamente las letras **y**, **n** o **c**, en cada uno de cuyos casos salvará los datos antes de terminar, terminará sin salvar, o volverá a la sesión de R. Los datos que se salvan estarán disponibles en la siguiente sesión de R.

Volver a trabajar con R es sencillo:

1. Haga que **trabajo** sea su directorio de trabajo e inicie el programa como antes:

```
$ cd trabajo
```

```
$ R
```

2. Dé las órdenes que estime convenientes a R y termine la sesión con la orden **q()**.

Bajo MACINTOSH o MICROSOFT WINDOWS el procedimiento a seguir es básicamente el mismo: Cree una carpeta o directorio. Ejecute R haciendo doble click en el icono correspondiente. Seleccione **New** dentro del menú **File** para indicar que desea iniciar un nuevo problema (lo que eliminará todos los objetos definidos dentro del espacio de trabajo) y a continuación seleccione **Save** dentro del menú **File** para salvar esta imagen en el directorio que acaba de crear. Puede comenzar ahora los análisis y cuando salga de R éste le preguntará si desea salvar la imagen en el directorio de trabajo.

Para continuar con este análisis posteriormente basta con pulsar en el icono de la imagen salvada, o bien puede ejecutar R y utilizar la opción **Open** dentro del menú **File** para seleccionar y abrir la imagen salvada.

## 1.5 Una sesión inicial

Se recomienda a los lectores que deseen ver cómo funciona R que realicen la sesión inicial dada en el Apéndice A, que se encuentra en la página 66.

## 1.6 Ayuda sobre funciones y capacidades

R contiene una ayuda similar a la orden `man` de UNIX. Para obtener información sobre una función concreta, por ejemplo `solve` la orden es

```
> help(solve)
```

Una forma alternativa es

```
> ?solve
```

Con las funciones especificadas por caracteres especiales, el argumento deberá ir entre comillas, para transformarlo en una "cadena de caracteres":

```
> help("[")
```

Podrá utilizar tanto comillas simples como dobles, y cada una de ellas puede utilizarse dentro de la otra, como en "Dijo 'Hola y adiós' y se marchó". En estas notas utilizaremos dobles comillas.

## 1.7 Órdenes de R. Mayúsculas y minúsculas

Técnicamente hablando, R es un *lenguaje de expresiones* con una sintaxis muy simple. Consecuentemente con sus orígenes en UNIX distingue entre mayúsculas y minúsculas, de tal modo que `A` y `a` son símbolos distintos y se referirán, por tanto, a variables distintas.

Las órdenes elementales consisten en *expresiones* o en *asignaciones*. Si una orden consiste en una expresión, se evalúa, se imprime y su valor se pierde. Una asignación, por el contrario, evalúa una expresión, no la imprime y guarda su valor en una variable.

Las órdenes se separan mediante punto y coma, (`;`), o por un cambio de línea. Si al terminar la línea la orden no está sintácticamente completa, R mostrará un nuevo signo de continuación, por ejemplo

```
+
```

en la línea siguiente y las sucesivas y continuará leyendo hasta que la orden esté sintácticamente completa. El signo de continuación puede ser modificado fácilmente. En estas notas omitiremos generalmente el símbolo de continuación y lo indicaremos mediante un sangrado.

## 1.8 Recuperación y corrección de órdenes previas

Bajo muchas versiones de UNIX, R permite recuperar y ejecutar órdenes previas. Las flechas verticales del teclado puede utilizarse para recorrer el *historial de órdenes*. Cuando haya recuperado una orden con este procedimiento, puede utilizar las flechas horizontales para desplazarse por él, puede eliminar caracteres con la tecla `Delete` o añadir nuevos caracteres.

La recuperación y edición pueden ser fácilmente adaptadas. Para ello debe buscar en el manual de UNIX la información sobre `readline`.

También puede utilizar el editor de textos `emacs` para trabajar más cómodamente de modo interactivo con R.

## 1.9 Ejecución de órdenes desde un archivo y redirección de la salida

Si tiene órdenes almacenadas en un archivo del sistema operativo, por ejemplo `órdenes.R` en el directorio de trabajo, `work`, puede ejecutarlas dentro de una sesión de R con la orden

```
> source("órdenes.R")
```



En las versiones de MACINTOSH y MICROSOFT WINDOWS, **Source** también está disponible dentro del menú **File** menu. Por otra parte, la orden **sink**, como por ejemplo en

```
> sink("resultado.lis")
```

enviará el resto de la salida, en vez de a la pantalla, al archivo del sistema operativo, **resultado.lis**, dentro del directorio de trabajo. La orden

```
> sink()
```

devuelve la salida de nuevo a la pantalla.

Si utiliza nombres absolutos de archivo en vez de nombre relativos, las órdenes se ejecutarán en ellos, independientemente del directorio de trabajo.

## 1.10 Almacenamiento y eliminación de objetos

Las entidades que R crea y manipula se denominan *objetos*. Estos pueden ser de muchos tipos: Variables, Variables indexadas numéricas, Cadenas de caracteres, funciones, etc. o incluso estructuras más complejas construidas a partir de otras más sencillas.

Durante una sesión de trabajo con R los objetos que se crean se almacenan por nombre (Discutiremos este proceso en la siguiente sección). La orden

```
> objects()
```

se puede utilizar para obtener los nombres de los objetos almacenados en R.

Para eliminar objetos puede utilizar la orden **rm**, por ejemplo:

```
> rm(x, y, z, tinta, chatarra, temporal, barra)
```

Los objetos creados durante una sesión de R pueden almacenarse en un archivo para su uso posterior. Al finalizar la sesión, R pregunta si desea hacerlo. En caso afirmativo todos los objetos se almacenan en el archivo **.RData**<sup>1</sup> en el directorio de trabajo.

En la siguiente ocasión que ejecute R cargará los objetos de este archivo así como el historial de órdenes.

Es recomendable que utilice un directorio de trabajo diferente para cada problema que analice con R. Es muy común crear objetos con los nombre **x** e **y**, por ejemplo. Estos nombres tienen sentido dentro de un análisis concreto, pero es muy difícil dilucidar su significado cuando se han realizado varios análisis en el mismo directorio.

---

<sup>1</sup>El punto inicial del nombre de este archivo indica que es *invisible* en UNIX.

## 2 Cálculos sencillos. Números y vectores

### 2.1 Vectores (numéricos). Asignación

R utiliza diferentes *estructuras de datos*. La estructura más simple es el *vector*, que es una colección ordenada de números. Para crear un vector, por ejemplo `x`, consistente en cinco números, por ejemplo 10.4, 5.6, 3.1, 6.4 y 21.7, use la orden

```
> x ← c(10.4, 5.6, 3.1, 6.4, 21.7)
```

Esta es una *asignación* en la que se utiliza la *función* `c()` que, en este contexto, puede tener un número arbitrario de vectores como *argumento* y cuyo valor es un vector obtenido mediante la concatenación de todos ellos.<sup>2</sup>

Un número, por sí mismo, se considera un vector de longitud uno.

Advierta que el operador de asignación, `(←)`, **no** es el operador habitual, `=`, que se reserva para otro propósito, sino que consiste en dos caracteres, `<` ('menor que') y `-` ('guión'), que obligatoriamente deben ir unidos y 'apuntan' al objeto que recibe el valor de la expresión. En este texto este operador se imprime "`←`", en vez de "`<-`" para facilitar la lectura.<sup>3</sup>

La asignación puede realizarse también mediante la función `assign()`. Una forma equivalente de realizar la asignación anterior es

```
> assign("x", c(10.4, 5.6, 3.1, 6.4, 21.7))
```

El operador usual, `←` ("`<-`"), puede interpretarse como una abreviatura de la función `assign()`.

Las asignaciones pueden realizarse también con una flecha apuntando a la derecha, realizando el cambio obvio en el operador de asignación. Por tanto, también podría escribirse

```
> c(10.4, 5.6, 3.1, 6.4, 21.7) -> x
```

Si una expresión se utiliza como una orden por sí misma, su valor se imprime y *se pierde*. Así pues, la orden

```
> 1/x
```

simplemente imprime los inversos de los cinco valores anteriores en la pantalla (por supuesto, el valor de `x` no se modifica).

Si a continuación hace la asignación

```
> y ← c(x, 0, x)
```

creará un vector, `y`, con 11 elementos, consistentes en dos copias de `x` con un cero entre ambas.

### 2.2 Aritmética vectorial

Los vectores pueden usarse en expresiones aritméticas, en cuyo caso las operaciones se realizan elemento a elemento. Dos vectores que se utilizan en la misma expresión no tienen por qué ser de la misma longitud. Si no lo son, el resultado será un vector de la longitud del más largo, y el más cortos será *reciclado* tantas veces como sea necesario (puede que no un número exacto de veces) hasta que coincida con el más largo. En particular, cualquier constante será simplemente repetida. De este modo, y siendo `x` e `y` los vectores antes definidos, la orden

```
> v ← 2*x + y + 1
```

<sup>2</sup>Con argumentos diferentes, por ejemplo *listas*, la acción de `c()` puede ser diferente. Véase §6.2.1.

<sup>3</sup>El símbolo de subrayado, "`_`", es un sinónimo del operador de asignación, pero no aconsejamos su utilización ya que produce un código menos legible.

genera un nuevo vector, `v`, de longitud 11, construido sumando, elemento a elemento, el vector `2*x` repetido 2.2 veces, el vector `y`, y el número 1 repetido 11 veces.

Los operadores aritméticos elementales son los habituales `+`, `-`, `*`, `/` y `^` para elevar a una potencia. Además están disponibles las funciones `log`, `exp`, `sin`, `cos`, `tan`, `sqrt`, bien conocidas.

Existen muchas más funciones, entre otras, las siguientes.

`max` y `min` que seleccionan respectivamente el mayor y el menor elemento de un vector, `range` cuyo valor es el vector de longitud dos `c(min(x), max(x))`, `length(x)` que es el número de elementos o longitud de `x`, `sum(x)` que es la suma de todos los elementos de `x`, y `prod(x)` que es el producto de todos ellos.

Dos funciones estadísticas son `mean(x)`, que calcula la media, esto es, `sum(x)/length(x)`, y `var(x)` que calcula la cuasi-varianza, esto es,

$$\text{sum}((x - \text{mean}(x))^2) / (\text{length}(x) - 1)$$

Si el argumento de `var()` es una matriz  $n \times p$ , el resultado es la matriz de cuasi-covarianzas  $p \times p$  correspondiente a interpretar las filas como vectores muestrales  $p$ -variantes.

Para ordenar un vector dispone de la función `sort(x)` que devuelve un vector del mismo tamaño que `x` con los elementos ordenados en orden creciente. También dispone de `order()` o `sort.list()` que produce la permutación que corresponde a la ordenación.

`rnorm(n)` es una función que genera un vector de números pseudoaleatorios de una distribución normal, del tamaño `n`.

## 2.3 Generación de sucesiones

En R existen varias funciones para generar sucesiones numéricas. Por ejemplo, `1:30` es el vector `c(1,2, ..., 29,30)`. El operador 'dos puntos' tiene máxima prioridad en una expresión, así, por ejemplo, `2*1:15` es el vector `c(2,4,6, ..., 28,30)`. Escriba `n <- 10` y compare las sucesiones `1:n-1` y `1:(n-1)`.

La forma `30:1` permite construir una sucesión descendente.

La función `seq()` permite generar sucesiones más complejas. Dispone de cinco argumentos, aunque no se utilizan todos simultáneamente. Si se dan los dos primeros indican el comienzo y el final de la sucesión, y si son los únicos argumentos, el resultado coincide con el operador 'dos puntos', esto es, `seq(2,10)` coincide con `2:10`.

Los argumentos de `seq()`, y de muchas funciones de R, pueden darse además por nombre, en cuyo caso el orden en que aparecen es irrelevante. En esta función los dos primeros argumentos se pueden dar por nombre mediante `from=valor-inicial` y `to=valor-final`; por tanto `seq(1,30)`, `seq(from=1, to=30)` y `seq(to=30, from=1)` son formas equivalente a `1:30`.

Los dos siguientes argumentos de `seq()` son `by=valor` y `length=valor`, y especifican el 'paso' y 'longitud' de la sucesión respectivamente. Si no se suministra ninguno, el valor predeterminado es `by=1` y `length` se calcula.

Por ejemplo

```
> seq(-5, 5, by=.2) -> s3
```

genera el vector `c(-5.0, -4.8, -4.6, ..., 4.6, 4.8, 5.0)` y lo almacena en `s3`. Similarmente

```
> s4 <- seq(length=51, from=-5, by=.2)
```

genera los mismos valores y los almacena en `s4`.

El quinto argumento de esta función es `along=vector`, y si se usa debe ser el único argumento, ya que crea una sucesión `1, 2, ..., length(vector)`, o la sucesión vacía si el vector es vacío (lo que puede ocurrir).

Una función relacionada con `seq` es `rep()` que permite duplicar un objeto de formas diversas. Su forma más sencilla es

```
> s5 ← rep(x, times=5)
```

que coloca cinco copias de `x`, una tras otra, y las almacena en `s5`.

## 2.4 Vectores lógicos

R no solo maneja vectores numéricos, sino también lógicos. Los elementos de un vector lógico solo pueden tomar dos valores: `FALSE` (falso) y `TRUE` (verdadero). Estos valores se representan también por `F` y `T`.

Los vectores lógicos aparecen al utilizar *condiciones*. Por ejemplo,

```
> temp ← x>13
```

almacena en `temp` un vector de la misma longitud de `x` y cuyos valores serán, respectivamente, `T` o `F` de acuerdo a que los elementos de `x` cumplan o no la condición `x>13`.

Los operadores lógicos son `<` (menor), `<=` (menor o igual), `>` (mayor), `>=` (mayor o igual), `==` (igual), y `!=` (distinto). Además, si `c1` y `c2` son expresiones lógicas, entonces `c1 & c2` es su intersección (“*conjunción*”), `c1 | c2` es su unión (“*disyunción*”) y `! c1` es la negación de `c1`.

Los vectores lógicos pueden utilizarse en expresiones aritméticas, en cuyo caso se transforman primero en vectores numéricos, de tal modo que `F` se transforma en `0` y `T` en `1`. Sin embargo hay casos en que un vector lógico y su correspondiente numérico no son equivalentes, como puede ver a continuación.

## 2.5 Valores faltantes

En ocasiones puede que no todas las componentes de un vector sean conocidas. Cuando falta un elemento, lo que se denomina ‘valor faltante’<sup>4</sup>, se le asigna un valor especial, `NA`<sup>5</sup>. En general, casi cualquier operación donde intervenga un `NA` da por resultado `NA`. La justificación es sencilla: Si no se puede especificar completamente la operación, el resultado no podrá ser conocido, y por tanto no estará disponible.

La función `is.na(x)` crea un vector lógico del tamaño de `x` cuyos elementos sólo valdrán `T` si el elemento correspondiente de `x` es `NA`, y `F` en caso contrario.

```
> ind ← is.na(z)
```

Nótese que la expresión lógica `x == NA` es distinta de `is.na(x)` puesto que `NA` no es realmente un valor, sino un indicador de una cantidad que no está disponible. Por tanto `x == NA` es un vector de la misma longitud de `x` con *todos* sus elementos `NA` puesto que la expresión lógica es incompleta.

Además hay una segunda clase de valores “perdidos”, producidos por el cálculo. Son los llamados valores `NaN`<sup>6</sup>. Este tipo de dato no existe en `S`, y por tanto se confunden con `NA` en `S-PLUS`. Ejemplos de `NaN` son

```
> 0/0
```

```
o
```

<sup>4</sup>En la literatura estadística inglesa, “missing value”

<sup>5</sup>Acrónimo en inglés de Not Available, no disponible.

<sup>6</sup>Acrónimo en inglés de “Not a Number”, esto es, “No es un número”.

```
> Inf - Inf
```

En resumen, `is.na(xx)` es TRUE *tanto* para los valores NA como para los NaN. Para diferenciar estos últimos existe la función `is.nan(xx)` que sólo es TRUE para valores NaN.

## 2.6 Vectores de caracteres

Las cadenas de caracteres o frases, también son utilizadas en R, por ejemplo, para etiquetar gráficos. Una cadena de caracteres se escribe entre comillas la sucesión de caracteres que la define, por ejemplo, "Altura" o "Resultados de la tercera iteración".

Los vectores de caracteres pueden concatenarse en un vector mediante la función `c()`.

Por otra parte, la función `paste()` une todos los vectores de caracteres que se le suministran y construye una sola cadena de caracteres. También admite argumentos numéricos, que convierte inmediatamente en cadenas de caracteres. En su forma predeterminada, en la cadena final, cada argumento original se separa del siguiente por un espacio en blanco, aunque ello puede cambiarse utilizando el argumento `sep="cadena"`, que sustituye el espacio en blanco por *cadena*, la cual podría ser incluso vacía.

Por ejemplo,

```
> labs ← paste(c("X","Y"), 1:10, sep=)
```

almacena, en `labs`, el vector de caracteres

```
("X1", "Y2", "X3", "Y4", "X5", "Y6", "X7", "Y8", "X9", "Y10")
```

Recuerde que al tener `c("X", "Y")` solo dos elementos, deberá repetirse 5 veces para obtener la longitud de 1:10.

## 2.7 Vectores de índices. Selección y modificación de subvectores

Puede seleccionar un subvector de un vector añadiendo al nombre del mismo un *vector de índices* entre corchetes, `[ y ]`. En general, de cualquier expresión cuyo resultado sea un vector, podrá obtener un subvector añadiéndole un vector de índices entre corchetes.

Los vectores de índices pueden ser de cuatro tipos distintos:

1. **Un vector lógico.** En este caso el vector de índices debe tener la misma longitud que el vector al que refiere. Sólo se seleccionarán los elementos correspondientes a valores T del vector de índices y se omitirá el resto. Por ejemplo,

```
> y ← x[!is.na(x)]
```

almacena en `y` los valores no-faltantes de `x`, en el mismo orden. Si `x` tiene valores faltantes, `y` será más corto que `x`. Análogamente,

```
> (x+1)[(!is.na(x)) & x>0] -> z
```

almacena en `z` los elementos del vector `x+1` para los que el correspondiente elemento de `x` es no-faltante y positivo.

2. **Un vector de números naturales positivos.** En este caso los elementos del vector de índices deben pertenecer al conjunto  $\{1, 2, \dots, \text{length}(x)\}$ . El resultado es un vector formado por los elementos del vector referido que corresponden a estos índices y en el orden en que aparecen en el vector de índices. El vector de índices puede tener cualquier longitud y el resultado será de esa misma longitud. Por ejemplo, `x[6]` es el sexto elemento de `x`, y

```
> x[1:10]
```

es el vector formado por los diez primeros elementos de `x`, (supuesto que `length(x) ≥ 10`).

Por otra parte,

```
> c("x","y")[rep(c(1,2,2,1), times=4)]
```

crea un vector de caracteres de longitud 16 formado por "x", "y", "y", "x" repetido cuatro veces.

- 3. Un vector de números naturales negativos.** En este caso, los índices indican los elementos del vector referido que deben *excluirse*. Así pues,

```
> y ← x[-(1:5)]
```

almacena en y todos los elementos de x excepto los cinco primeros (suponiendo que x tiene al menos cinco elementos).

- 4. Un vector de caracteres.** Esta opción solo puede realizarse si el vector posee el atributo `names` (nombres) para identificar sus componentes, en cuyo caso se comportará de modo similar al punto 2.

```
> fruta ← c(5, 10, 1, 20)
```

```
> names(fruta) ← c("naranja", "plátano", "manzana", "pera")
```

```
> postre ← fruta[c("manzana","naranja")]
```

La ventaja en este caso es que los *nombres* son a menudo más fáciles de recordar que los *índices numéricos*. Esta opción es especialmente útil al tratar de la estructura de “hoja de datos” (data frame) que veremos posteriormente.

La variable de almacenamiento también puede ser indexada, en cuyo caso la asignación se realiza solamente sobre los elementos referidos.

El vector asignado debe ser de la misma longitud del vector de índices y, en el caso de un vector de índices lógico, debe ser de la misma longitud del vector que indexa. Por ejemplo,

```
> x[is.na(x)] ← 0
```

sustituye cada valor faltante de x por un cero. Por otra parte,

```
> y[y<0] ← -y[y<0]
```

equivale a

```
> y ← abs(y)7
```

---

<sup>7</sup>Tenga en cuenta que `abs()` no se comporta correctamente con números complejos, en ellos debería usar `Mod()`.

## 3 Objetos: Modos y atributos

### 3.1 Atributos intrínsecos: *modo* y *longitud*

Las entidades que manipula R se conocen con el nombre de *objetos*. Por ejemplo, los vectores de números, reales o complejos, los vectores lógicos o los vectores de caracteres. Este tipo de objetos se denominan estructuras ‘atómicas’ puesto que todos sus elementos son del mismo tipo o *modo*, bien sea *numeric*<sup>8</sup> (numérico), *complex* (complejo), *logical* (lógico) o *character* (carácter).

Los elementos de un vector deben ser *todos del mismo modo* y éste será el modo del vector. Esto es, un vector será, completo, de modo *logical*, *numeric*, *complex* o *character*. La única excepción a esta regla es que cualquiera de ellos puede contener el valor NA. Debe tener en cuenta que un vector puede ser vacío, pero pese a ello tendrá un modo. Así, el vector de caracteres vacío aparece como `character(0)` y el vector numérico vacío aparece como `numeric(0)`.

R también maneja objetos denominados *listas* que son del modo *list* (lista) y que consisten en sucesiones de objetos, cada uno de los cuales puede ser de un modo distinto. Las *listas* se denominan estructuras ‘recursivas’ puesto que sus componentes pueden ser asimismo listas.

Existen otras estructuras recursivas que corresponden a los modos *function* (función) y *expression* (expresión). El modo *function* está formado por las funciones que constituyen R, unidas a las funciones escritas por cada usuario, y que discutiremos más adelante. El modo *expression* corresponde a una parte avanzada de R que no trataremos aquí, excepto en lo mínimo necesario para el tratamiento de *fórmulas* en la descripción de modelos estadísticos.

Con el *modo* de un objeto designamos el tipo básico de sus constituyentes fundamentales. Es un caso especial de un *atributo* de un objeto. Los *atributos* de un objeto suministran información específica sobre el propio objeto. Otro atributo de un objeto es su *longitud*. Las funciones `mode(objeto)` y `length(objeto)` se pueden utilizar para obtener el modo y longitud de cualquier estructura.

Por ejemplo, si `z` es un vector complejo de longitud 100, entonces `mode(z)` es la cadena "complex", y `length(z)` es 100.

R realiza cambios de modo cada vez que se le indica o es necesario (y también en algunas ocasiones en que no lo es). Por ejemplo, si escribe

```
> z ← 0:9
```

y a continuación escribe

```
> digitos ← as.character(z)
```

el vector `digitos` será el vector de caracteres ("0", "1", "2", ..., "9"). Si a continuación aplica un nuevo cambio de modo

```
> d ← as.integer(digitos)
```

R reconstruirá el vector numérico de nuevo y, en este caso, `d` y `z` coinciden.<sup>9</sup> Existe una colección completa de funciones de la forma `as.lo-que-sea()`, tanto para forzar el cambio de modo, como para asignar un atributo a un objeto que carece de él. Es aconsejable consultar la ayuda para familiarizarse con estas funciones.

### 3.2 Modificación de la longitud de un objeto

Un objeto, aunque esté “vacío”, tiene modo. Por ejemplo,

<sup>8</sup>El modo *numérico* consiste realmente en dos modos distintos, *integer* (entero) y *double* (doble precisión) como se explica en el manual.

<sup>9</sup>En general, al forzar el cambio de numérico a carácter y de nuevo a numérico, no se obtienen los mismos resultados, debido, entre otros, a los errores de redondeo.

```
> v ← numeric()
```

almacena en `v` una estructura vacía de vector numérico. Del mismo modo, `character()` es un vector de caracteres vacío, y lo mismo ocurre con otros tipos. Una vez creado un objeto con un tamaño cualquiera, pueden añadirse nuevos elementos sin más que asignarlos a un índice que esté fuera del rango previo. Por ejemplo,

```
> v[3] ← 17
```

transforma `v` en un vector de longitud 3, (cuyas dos primeras componentes serán `NA`). Esta regla se aplica a cualquier estructura, siempre que los nuevos elementos sean compatibles con el modo inicial de la estructura.

Este ajuste automático de la longitud de un objeto se utiliza a menudo, por ejemplo en la función `scan()` para entrada de datos. (Véase §7.2.)

Análogamente, puede reducirse la longitud de un objeto sin más que realizar una nueva asignación. Si, por ejemplo, `alfa` es un objeto de longitud 10, entonces

```
> alfa ← alfa[2 * 1:5]
```

lo transforma en un objeto de longitud 5 formado por los elementos de posición par del objeto inicial.

### 3.3 `attributes()` y `attr()`

La función `attributes(objeto)` proporciona una lista de todos los atributos no intrínsecos definidos para el objeto en ese momento. La función `attr(objeto, nombre)` puede usarse para seleccionar un atributo específico. Estas funciones no se utilizan habitualmente, sino que se reservan para la creación de un nuevo atributo con fines específicos, por ejemplo para asociar una fecha de creación o un operador con un objeto de R. Sin embargo, es un concepto muy importante que no debe olvidarse.

La asignación o eliminación de atributos de un objeto debe realizarse con precaución, ya que los atributos forman parte del sistema de objetos utilizados en R.

Cuando se utiliza en la parte que recibe la asignación, puede usarse para asociar un nuevo atributo al `objeto` o para cambiar uno existente. Por ejemplo,

```
> attr(z, "dim") ← c(10, 10)
```

permite a R tratar `z` como si fuese una matriz de  $10 \times 10$ .

### 3.4 Clases de objetos

Cada objeto pertenece a una *clase*, y ello permite utilizar en R programación dirigida a objetos.

Por ejemplo, si un objeto pertenece a la clase `data.frame`, se imprimirá de un modo especial; cuando le aplique la función `plot()` ésta mostrará un gráfico de un tipo especial; y otras funciones genéricas, como `summary()`, producirán un resultado especial; todo ello en función de la pertenencia a dicha clase.

Para eliminar temporalmente los efectos de la clase puede utilizar la función `unclass()`. Por ejemplo, si `invierno` pertenece a la clase `data.frame`, entonces

```
> invierno
```

escribe el objeto en la forma de la clase, en tanto que

```
> unclass(invierno)
```



lo imprime como una lista ordinaria. Sólo debe utilizar esta función en situaciones muy concretas, como por ejemplo, cuando haga pruebas para comprender el concepto de clase y de función genérica.

Las clases y las funciones genéricas serán tratadas muy brevemente en §9.8.

## 4 Factores Nominales y Ordinales

Un *factor* es un vector utilizado para especificar una clasificación discreta de los elementos de otro vector de igual longitud. En R existen factores *nominales* y factores *ordinales*.

### 4.1 Un ejemplo específico

Suponga que dispone de una muestra de 30 personas de Australia<sup>10</sup> de tal modo que su estado o territorio se especifica mediante un vector de caracteres con las abreviaturas de los mismos:

```
> estado <- c("tas", "sa", "qld", "nsw", "nsw", "nt", "wa", "wa",
              "qld", "vic", "nsw", "vic", "qld", "qld", "sa", "tas",
              "sa", "nt", "wa", "vic", "qld", "nsw", "nsw", "wa",
              "sa", "act", "nsw", "vic", "vic", "act")
```

Recuerde que, para un vector de caracteres, la palabra “ordenado” indica que está en orden alfabético.

Un *factor* se crea utilizando la función `factor()`:

```
> FactorEstado <- factor(estado)
```

La función `print()` trata a los factores de un modo distinto que a los vectores ordinarios:

```
> FactorEstado
[1] tas sa qld nsw nsw nt wa wa qld vic nsw vic qld qld sa
[16] tas sa nt wa vic qld nsw nsw wa sa act nsw vic vic act
```

Puede utilizar la función `levels()` para ver los niveles de un factor:

```
> levels(FactorEstado) [1] "act" "nsw" "nt" "qld" "sa" "tas" "vic" "wa"
```

### 4.2 La función `tapply()`. Variables desastradas (ragged arrays)

Como continuación del ejemplo anterior, suponga que disponemos en otro vector de los ingresos de las mismas personas (en unidades apropiadas)

```
> ingresos <- c(60, 49, 40, 61, 64, 60, 59, 54, 62, 69, 70, 42, 56,
                61, 61, 61, 58, 51, 48, 65, 49, 49, 41, 48, 52, 46,
                59, 46, 58, 43)
```

Para calcular la media muestral para cada estado podemos usar la función `tapply()`:

```
> MediaIngresos <- tapply(ingresos, FactorEstado, mean)
```

que devuelve el vector de medias con las componentes etiquetadas con los niveles:

```
> MediaIngresos
      act      nsw      nt      qld      sa      tas      vic      wa
44.50000 57.33333 55.50000 53.60000 55.00000 60.50000 56.00000 52.25000
```

La función `tapply()` aplica una función, en este ejemplo la función `mean()`, a cada grupo de componentes del primer argumento, en este ejemplo `ingresos`, definidos por los niveles del segundo argumento, en este ejemplo `FactorEstado`, como si cada grupo fuese un vector por sí

<sup>10</sup>Para quienes no conocen la estructura administrativa de Australia, existen ocho estados y territorios en la misma: Australian Capital Territory, New South Wales, Northern Territory, Queensland, South Australia, Tasmania, Victoria, y Western Australia; y sus correspondientes abreviaturas son: act, nsw, nt, qld, sa, tas, vic, y wa.

solo. El resultado es una estructura de la misma longitud que niveles posee el factor. Puede consultar la ayuda para obtener más detalles.

Suponga que ahora desea calcular las desviaciones típicas de las medias de ingresos por estados. Para ello es necesario escribir una función en R que calcule la desviación típica de un vector. Aunque aún no se ha explicado en este texto cómo escribir funciones<sup>11</sup>, puede admitir que existe la función `var()` que calcula la varianza muestral o cuasi-varianza, y que la función buscada puede construirse con la asignación:

```
> StdErr ← function(x) sqrt(var(x)/length(x))
```

Ahora puede calcular los valores buscados mediante

```
> ErrorTipicoIngresos ← tapply(ingresos, FactorEstado, StdErr)
```

con el siguiente resultado:

```
> ErrorTipicoIngresos
      act      nsw      nt      qld      sa      tas      vic      wa
1.500000 4.310195 4.500000 4.106093 2.738613 0.500000 5.244044 2.657536
```

Como ejercicio puede calcular el intervalo de confianza al 95% de la media de ingresos por estados. Para ello puede utilizar la función `tapply()`, la función `length()` para calcular los tamaños muestrales, y la función `qt()` para encontrar los percentiles de las distribuciones  $t$  de Student correspondientes.

La función `tapply()` puede utilizarse para aplicar una función a un vector indexado por diferentes categorías simultáneamente. Por ejemplo, para dividir la muestra tanto por el estado como por el sexo. Los elementos del vector se dividirán en grupos correspondientes a las distintas categorías y se aplicará la función a cada uno de dichos grupos. El resultado es una variable indexada etiquetada con los niveles de cada categoría.

La combinación de un vector<sup>12</sup> y un factor para etiquetarlo es un ejemplo de lo que se llama *variable indexada desastrada* (ragged array) puesto que los tamaños de las subclases son posiblemente irregulares. Cuando estos tamaños son iguales la indexación puede hacerse implícitamente y además más eficientemente, como veremos a continuación.

---

<sup>11</sup>La escritura de funciones será tratada en §9.

<sup>12</sup>En general de una variable indexada

## 5 Variables indexadas. Matrices

### 5.1 Variables indexadas (Arrays)

Una variable indexada (array) es una colección de datos, por ejemplo numéricos, indexada por varios índices. R permite crear y manipular variables indexadas en general y en particular, matrices.

Un vector de dimensiones es un vector de números enteros positivos. Si su longitud es  $k$  entonces la variable indexada correspondiente es  $k$ -dimensional. Los elementos del vector de dimensiones indican los límites superiores de los  $k$  índices. Los límites inferiores siempre valen 1.

Un vector puede transformarse en una variable indexada cuando se asigna un vector de dimensiones al atributo *dim*. Supongamos, por ejemplo, que **z** es un vector de 1500 elementos. La asignación

```
> dim(z) ← c(3,5,100)
```

hace que R considere a **z** como una variable indexada de dimensión  $3 \times 5 \times 100$ .

Existen otras funciones, como `matrix()` y `array()`, que permiten asignaciones más sencillas y naturales, como se verá en §5.4.

Los elementos del vector pasan a formar parte de la variable indexada siguiendo la regla<sup>13</sup> de que el primer índice es el que se mueve más rápido y el último es el más lento.

Por ejemplo, si se define una variable indexada **a**, con vector de dimensiones `c(3,4,2)`, la variable indexada tendrá  $3 \times 4 \times 2 = 24$  elementos que se formarán a partir de los elementos originales en el orden `a[1,1,1]`, `a[2,1,1]`, ..., `a[2,4,2]`, `a[3,4,2]`.

### 5.2 Elementos de una variable indexada

Un elemento de una variable indexada puede referirse dando el nombre de la variable y, entre corchetes, los índices que lo refieren, separados por comas.

En general, puede referir una parte de una variable indexada mediante una sucesión de *vectores índices*, teniendo en cuenta que *si un vector índice es vacío, equivale a utilizar todo el rango de valores para dicho índice*.

Así, en el ejemplo anterior, `a[2,,]` es una variable indexada  $4 \times 2$ , con vector de dimensión `c(4,2)` y sus elementos son

```
a[2,1,1], a[2,2,1], a[2,3,1], a[2,4,1], a[2,1,2], a[2,2,2], a[2,3,2], a[2,4,2],
```

en ese orden. A su vez, `a[, ,]` equivale a la variable completa, que coincide con omitir completamente los índices y utilizar simplemente **a**.

Para cualquier variable indexada, por ejemplo **Z**, el vector de dimensión puede referirse explícitamente mediante `dim(Z)` (en cualquiera de las dos partes de una asignación).

Asimismo, si especifica una variable indexada con *un solo índice o vector índice*, sólo se utilizan los elementos correspondientes del vector de datos, y el vector de dimensión se ignora. En caso de que el índice no sea un vector, sino a su vez una variable indexada, el tratamiento es distinto, como ahora veremos.

---

<sup>13</sup>Esta es la misma que se utiliza en el lenguaje Fortran

### 5.3 Uso de variables indexadas como índices

Una variable indexada puede utilizar no sólo un vector de índices, sino incluso una *variable indexada de índices* tanto para asignar un vector a una colección irregular de elementos de una variable indexada o para extraer una colección irregular de elementos.

Veamos un ejemplo sobre una matriz, que es una variable indexada con dos índices. Puede construirse un índice matricial consistente en dos columnas y varias filas. Los elementos del índice matricial son los índices fila y columna para construir la matriz de índices. Supongamos que **X** es una variable indexada  $4 \times 5$  y que desea hacer lo siguiente:

- Extraer los elementos **X**[1,3], **X**[2,2] y **X**[3,1] con una estructura de vector, y
- Reemplazar dichos elementos de **X** con ceros.

Para ello puede utilizar una matriz de índices de  $3 \times 2$  como en el ejemplo de la figura 1.

---

```
> x ← array(1:20,dim=c(4,5)) # Genera una variable indexada 4 × 5.
> x
      [,1] [,2] [,3] [,4] [,5]
[1,]    1    5    9   13   17
[2,]    2    6   10   14   18
[3,]    3    7   11   15   19
[4,]    4    8   12   16   20
> i ← array(c(1:3,3:1),dim=c(3,2))
> i
      [,1] [,2]
[1,]    1    3
[2,]    2    2
[3,]    3    1
> x[i]
      # Extrae los elementos.
[1] 9 6 3
> x[i] ← 0
      # Sustituye los elementos por ceros.
> x
      [,1] [,2] [,3] [,4] [,5]
[1,]    1    5    0   13   17
[2,]    2    0   10   14   18
[3,]    0    7   11   15   19
[4,]    4    8   12   16   20
>
```

---

Figura 1: Uso de una matriz de índices

---

Un ejemplo algo más complejo consiste en generar la matriz de diseño de un diseño en bloques definido por los factores **bloques** (niveles **b**) y **variedades** (niveles **v**), y que el número de parcelas es **n**. Puede hacerlo del siguiente modo:

```
> Xb ← matrix(0, n, b)
> Xv ← matrix(0, n, v)
> ib ← cbind(1:n, bloques)
> iv ← cbind(1:n, variedades)
> Xb[ib] ← 1
> Xv[iv] ← 1
> X ← cbind(Xb, Xv)
```

Además, puede construir la matriz de incidencia, **N**, mediante

```
> N ← crossprod(Xb, Xv)
```

También puede construirla directamente mediante la función `table()`:

```
> N ← table(bloques, variedades)
```

## 5.4 La función `array()`

Una variable indexada puede construirse directamente a partir de un vector mediante la función `array`, que tiene la forma

```
> Z ← array(vector_de_datos, vector_de_dimensiones)
```

Por ejemplo, si el vector `h` contiene 24 números (o menos), la orden

```
> Z ← array(h, dim=c(3,4,2))
```

usa `h` para almacenar en `Z` una variable indexada de dimensión  $3 \times 4 \times 2$ . Si el tamaño de `h` es exactamente 24, el resultado coincide con el de

```
> dim(Z) ← c(3,4,2)
```

Sin embargo, si `h` es más corto de 24, sus valores se repiten desde el principio tantas veces como sea necesario para obtener 24 elementos. Puede ver §5.4.1. El caso extremo (muy común) corresponde a un vector de longitud 1, como en el siguiente ejemplo

```
> Z ← array(0, c(3,4,2))
```

en que `Z` es una variable indexada compuesta enteramente de ceros.

Además, `dim(Z)`, el vector de dimensiones, es el vector `c(3,4,2)`, `Z[1:24]`, es un vector de datos que coincide con `h`, y tanto `Z[]`, con índice vacío, como `Z`, sin índices, son la variable indexada con estructura de variable indexada.

Las variables indexadas pueden utilizarse en expresiones aritméticas y el resultado es una variable indexada formada a partir de las operaciones elemento a elemento de los vectores subyacentes. Los atributos `dim` de los operandos deben ser iguales en general y coincidirá con el vector de dimensiones del resultado. Así pues, si `A`, `B` y `C` son variables indexadas similares, entonces

```
> D ← 2*A*B + C + 1
```

almacena en `D` una variable indexada similar, cuyo vector de datos es el resultado de las operaciones indicadas sobre los vectores de datos subyacentes a `A`, `B` y `C`. Las reglas exactas correspondientes a cálculos en que se mezclan variables indexadas y vectores se estudian a continuación.

### 5.4.1 Operaciones con variables indexadas y vectores. Reciclado.

Cuando se realizan operaciones que mezclan variables indexadas y vectores, se siguen los siguientes criterios:

- La expresión se analiza de izquierda a derecha.
- Si un vector es más corto, se extiende repitiendo sus elementos (lo que se denomina reciclado) hasta alcanzar el tamaño del vector más largo.
- Si sólo hay variables indexadas y vectores más cortos, las variables indexadas deben tener el mismo atributo `dim` o se producirá un error.
- Si hay un vector más largo que una variable indexada anterior, se produce el siguiente mensaje de error:

```
Error: dim<- length of dims do not match the length of object
```

- Si hay variables indexadas y no se produce error, el resultado es una variable indexada del mismo atributo `dim` que las variables indexadas que intervienen en la operación.

## 5.5 Producto exterior de dos variables indexadas

Una operación de importancia fundamental entre variables indexadas es el *producto exterior*. Si **a** y **b** son dos variables indexadas numéricas, su producto exterior es una variable indexada cuyo vector de dimensión es la concatenación de los correspondientes de los operandos, en el orden de la operación, y cuyo vector de datos subyacente se obtiene mediante todos los posibles productos de los elementos de los vectores subyacentes de **a** y **b**. El producto exterior se obtiene mediante el operador `%o%`:

```
> ab ← a %o% b
```

o bien, en forma funcional, mediante `outer`:

```
> ab ← outer(a, b, '*')
```

La función “multiplicación” puede reemplazarse por cualquier función de dos variables. Por ejemplo, para calcular la función

$$f(x, y) = \frac{\cos(y)}{1 + x^2}$$

sobre la retícula formada por todos los puntos obtenidos combinando las ordenadas y abscisas definidas por los elementos de los vectores **x** e **y** respectivamente, puede utilizar<sup>14</sup> las órdenes:

```
> f ← function(x,y) cos(y)/(1 + x^2)
> z ← outer(x, y, f)
```

En particular, el producto exterior de dos vectores, es una variable indexada con dos índices (esto es, una matriz, de rango 1 a lo sumo). Debe tener en cuenta que el producto exterior no es conmutativo.

### 5.5.1 Ejemplo: Distribución del determinante de una matriz $2 \times 2$ de dígitos

Un ejemplo apropiado se presenta en el cálculo del determinante de una matriz  $2 \times 2$ ,  $\begin{bmatrix} a & b \\ c & d \end{bmatrix}$ , en que cada elemento de la misma es un número natural entre 0 y 9. El problema planteado consiste en encontrar la distribución de los determinantes,  $ad - bc$ , de todas las matrices posibles de esta forma y representarla gráficamente, supuesto que cada dígito se elige al azar de una distribución uniforme.

Para ello puede utilizar la función `outer()` dos veces:

```
> d ← outer(0:9, 0:9)
> fr ← table(outer(d, d, "-"))
> plot(as.numeric(names(fr)), fr, type="h",
       xlab="Determinante", ylab="Frecuencia")
```

Advierta como se ha forzado el atributo `names` de la tabla de frecuencias a numérico, para recuperar el rango de los valores de los determinantes. La forma aparentemente “obvia” de resolver este problema mediante iteraciones de tipo `for`, que se discutirán en §8.2, es tan ineficiente que es impracticable. Tal vez le haya sorprendido que aproximadamente una de cada veinte matrices sea singular.

<sup>14</sup>La definición de una función en R se estudia en el capítulo 9.

## 5.6 Traspuesta generalizada de una variable indexada

La función `aperm(a, perm)` puede usarse para permutar la variable indexada `a`. El argumento `perm` debe ser una permutación de los enteros  $\{1, 2, \dots, k\}$ , siendo  $k$  el número de índices de `a`. El resultado es una variable indexada del mismo tamaño que `a` en la que la dimensión que en la original era `perm[j]` será ahora la dimensión `j`. Si `A` es una matriz

```
> B ← aperm(A, c(2,1))
```

entonces `B` es la matriz traspuesta de `A`. En el caso de matrices es más sencillo utilizar la función `t()`, y bastaría escribir `B ← t(A)`.

## 5.7 Operaciones con matrices: Producto matricial. Inversa de una matriz. Resolución de sistemas lineales

Como ya se ha indicado varias veces, una matriz es simplemente una variable indexada con dos índices. Ahora bien, su importancia es tal que necesita un apartado especial. R dispone de muchos operadores y funciones diseñados específicamente para matrices. Por ejemplo, acabamos de ver que `t(X)` es la matriz traspuesta de `X`. Las funciones `nrow` y `ncol` devuelven el número de filas y de columnas de una matriz.

El operador `%*%` realiza el producto matricial. Una matriz de  $n \times 1$  o de  $1 \times n$  puede ser utilizada como un vector  $n$ -dimensional en caso necesario. Análogamente R puede usar automáticamente un vector en una operación matricial convirtiéndolo en una matriz fila o una matriz columna cuando ello es posible (A veces la conversión no puede hacerse de modo único, como veremos después).

Si, por ejemplo, `A` y `B`, son matrices cuadradas del mismo tamaño, entonces

```
> A * B
```

es la matriz de productos elemento a elemento, en tanto que

```
> A %*% B
```

es el producto matricial. Si `x` es un vector (de la dimensión apropiada) entonces

```
> x %*% A %*% x
```

es una forma cuadrática.<sup>15</sup>

La función `crossprod()` realiza el producto cruzado de matrices, esto es

```
> crossprod(X, y) coincide con t(X) %*% y
```

pero la operación es más eficiente. Si omite el segundo argumento de la función `crossprod()`, ésta lo toma igual al primero.

Otra función matricial importante es `solve(A, b)` que resuelve el sistema de ecuaciones  $Ax = b$ . Si omite el segundo argumento, `solve(A)`, toma  $b = I$ , y obtiene la matriz inversa de `A`. Además existen `svd()` que calcula los autovalores y autovectores de una matriz rectangular y su descomposición SVD, `qr()` que realiza la descomposición QR, y `eigen()` que calcula los autovalores y autovectores de una matriz simétrica.

También existe la función `diag()`. Si su argumento es una matriz, `diag(matriz)`, devuelve un vector formado por los elementos de la diagonal de la misma. Si, por el contrario, su argumento es un vector (de longitud mayor que uno), `diag(vector)`, lo transforma en una matriz diagonal

<sup>15</sup>Si hubiese escrito `x %*% x` el resultado es ambiguo, pues tanto podría significar  $x'x$  como  $xx'$ , donde  $x$  es la forma columna. En este tipo de casos, la interpretación corresponde a la matriz de menor tamaño, por lo que en este ejemplo el resultado es el escalar  $x'x$ . La matriz  $xx'$  puede calcularse, bien mediante `cbind(x) %*% x`, mediante `x %*% rbind(x)` o mediante `x %*% rbind(x)`, puesto que tanto el resultado de `rbind()` como el de `cbind()` son matrices.



cuyos elementos diagonales son los del vector. Y, por último, si su argumento es un número natural, `n`, lo transforma en una matriz identidad de tamaño  $n \times n$ .

En esta colección de funciones matriciales, sorprendentemente, falta una que calcule el determinante de una matriz cuadrada, aunque éste puede calcularse fácilmente como producto de los autovalores, como veremos posteriormente.

## 5.8 Submatrices. Funciones `cbind()` y `rbind()`.

Las funciones `cbind()` y `rbind()` construyen matrices uniendo otras matrices (o vectores), horizontalmente (modo columna) o verticalmente (modo fila), respectivamente.

En la asignación

```
> X ← cbind(arg1, arg2, arg3, ...)
```

los argumentos pueden ser vectores de cualquier longitud o matrices con el mismo número de filas. El resultado es una matriz cuyas columnas son los argumentos concatenados, `arg1`, `arg2`, ....

Si alguno de los argumentos de `cbind()` es un vector, y hay alguna matriz, el vector no puede ser más largo que el número de filas de las matrices presentes, y si es más corto, se recicla hasta alcanzar el número indicado. Si sólo hay vectores, los más cortos se reciclan hasta alcanzar el tamaño del mayor.

La función `rbind()` realiza el mismo papel, sustituyendo filas por columnas.

Supongamos que `X1` y `X2` tienen el mismo número de filas. Para combinar las columnas de ambas en una matriz `X`, que tendrá el mismo número de filas, y añadirle una columna inicial de unos, puede escribir

```
> X ← cbind(1, X1, X2)
```

El resultado de `cbind()` o de `rbind()` siempre es una matriz y estas funciones constituyen, por tanto, la forma más sencilla para tratar un vector como una matriz columna o una matriz fila, respectivamente.

## 5.9 La función de concatenación, `c()`, con variables indexadas

En tanto que `cbind()` y `rbind()` son funciones de concatenación que respetan el atributo `dim`, la función `c()` no lo hace, sino que despoja a los objetos numéricos de los atributos `dim` y `dimnames`, lo que es útil en determinadas situaciones.

La forma *oficial* de transformar una variable indexada en el vector subyacente es utilizar la función `as.vector()`,

```
> vec ← as.vector(X)
```

Sin embargo, se obtiene un resultado análogo utilizando la función `c()` debido al efecto colateral citado:

```
> vec ← c(X)
```

Existen sutiles diferencias entre ambos resultados, aunque fundamentalmente la elección entre ambas es una cuestión de estilo(aunque personalmente preferimos la primera forma).

## 5.10 Tablas de frecuencias a partir de factores. La función `table()`

Recuerde que un factor define una partición en grupos. Del mismo modo, dos factores definen una tabla de doble entrada, y así sucesivamente. La función `table()` calcula tablas a partir de

factores de igual longitud. Si existen  $k$  argumentos, el resultado será una variable de  $k$  entradas, de frecuencias.

Recuerde que en un ejemplo anterior, `FactorEstado` era un factor que indicaba el estado de procedencia. La asignación

```
> FrecEstado ← table(FactorEstado)
```

almacena en `FrecEstado` una tabla de las frecuencias de cada estado en la muestra. Las frecuencias se ordenan y etiquetan con los niveles del factor. Esta orden es equivalente, y más sencilla, que

```
> FrecEstado ← tapply(FactorEstado, FactorEstado, length)
```

Suponga ahora que `FactorIngresos` es un factor que define “tipos de ingresos”, por ejemplo, mediante la función `cut()`:

```
> factor(cut(ingresos,breaks=35+10*(0:7))) -> FactorIngresos
```

Entonces, puede calcular una tabla de frecuencias de doble entrada del siguiente modo:

```
> table(FactorIngresos,FactorEstado)
      FactorEstado
FactorIngresos act nsw nt qld sa tas vic wa
(35,45]      1   1  0   1  0  0   1  0
(45,55]      1   1  1   1  2  0   1  3
(55,65]      0   3  1   3  2  2   2  1
(65,75]      0   1  0   0  0  0   1  0
```

La extensión a tablas de frecuencias de varias entradas es inmediata.

## 6 Listas y hojas de datos. Utilización

### 6.1 Listas

Una *lista* es un objeto consistente en una colección ordenada de objetos, conocidos como *componentes*.

No es necesario que los componentes sean del mismo modo, así una lista puede estar compuesta de, por ejemplo, un vector numérico, un valor lógico, una matriz y una función. El siguiente es un ejemplo de una lista:

```
> Lst ← list(nombre="Pedro", esposa="María", no.hijos=3, edad.hijos=c(4,7,9))
```

Los componentes siempre están *numerados* y pueden ser referidos por dicho número. En este ejemplo, `Lst` es el nombre de una lista con cuatro componentes, cada uno de los cuales puede ser referido, respectivamente, por `Lst[[1]]`, `Lst[[2]]`, `Lst[[3]]` y `Lst[[4]]`. Como, además, `Lst[[4]]` es un vector, `Lst[[4]][1]` refiere su primer elemento.

La función `length()` aplicada a una lista devuelve el número de componentes (del primer nivel) de la lista.

Los componentes de una lista pueden tener *nombre*, en cuyo caso pueden ser referidos también por dicho nombre, mediante una expresión de la forma

```
> nombre_de_lista$nombre_de_componente
```

Esta es una convención útil que facilita la obtención de una componente cuando se ha olvidado su número.

En el ejemplo anterior,

`Lst$nombre` coincide con `Lst[[1]]` y vale "Pedro",

`Lst$esposa` coincide con `Lst[[2]]` y vale "María",

`Lst$edad.hijos[1]` coincide con `Lst[[4]][1]` y vale 4.

También es posible utilizar los nombres de los componentes entre dobles corchetes, por ejemplo, `Lst[["nombre"]]` coincide con `Lst$nombre`. Esta opción es muy útil en el caso en que el nombre de los componentes se almacena en otra variable, como en

```
> x ← "nombre"; Lst[[x]]
```

Es muy importante distinguir claramente entre `Lst[[1]]` y `Lst[1]`. “[...]” es el operador utilizado para seleccionar un sólo elemento, mientras que “[...]” es un operador general de indexado. Esto es, `Lst[[1]]` es el *primer objeto de la lista* `Lst`, y si es una lista con nombres, el nombre *no* está incluido. Por su parte, `Lst[1]`, es una *sublista de la lista* `Lst` consistente en la primera componente. Si la lista tiene nombre, éste se transfiere a la sublista.

Los nombres de los componentes pueden abreviarse hasta el mínimo de letras necesarios para identificarlos de modo exacto. Así, en

```
> Lista ← list(coeficientes=c(1.3,4), covarianza=.87)
```

`Lst$coeficientes` puede especificarse mediante `Lista$coe`, y `Lista$covarianza` como `Lista$cov`.

El vector de nombres es un atributo de la lista, y como el resto de atributos puede ser manipulado. Otras estructuras, además de las listas, también pueden tener el atributo *names*.

## 6.2 Construcción y modificación de listas

La función `list()` permite crear listas a partir de objetos ya existentes. Una asignación de la forma

```
> Lista ← list(nombre1=objeto1, ..., nombrem=objetom)
```

almacena en `Lista` una lista de  $m$  componentes que son `objeto1`, ..., `objetom`, a los cuales asigna los nombres `nombre1`, ..., `nombrem`, que pueden ser libremente elegidos<sup>16</sup>. Si omite los nombres, las componentes sólo estarán numeradas. Las componentes se *copian* para construir la lista y los originales no se modifican.

Las listas, como todos los objetos indexados, pueden ampliarse especificando componentes adicionales. Por ejemplo

```
> Lst[5] ← list(matriz=Mat)
```

### 6.2.1 Concatenación de listas

Al suministrar listas como argumentos a la función `c()` el resultado es una lista, cuyos componentes son todos los de los argumentos unidos sucesivamente.

```
> lista.ABC ← c(lista.A, lista.B, lista.C)
```

Recuerde que cuando los argumentos eran vectores, esta función los unía todos en un único vector. En este caso, el resto de atributos, como `dim`, se pierden.

## 6.3 Algunas funciones que devuelven una lista

Las funciones y expresiones de R sólo devuelven un objeto como resultado, por tanto, si deben devolver varios objetos, previsiblemente de diferentes tipos, la forma usual es una lista con nombres.

### 6.3.1 Autovalores y autovectores

Como ya hemos indicado, la función `eigen()` calcula los autovalores y autovectores de una matriz simétrica. El resultado es una lista de dos componentes llamados `values` y `vectors`. La asignación

```
> ev ← eigen(Sm)
```

almacenará esta lista en `ev`. Por tanto `ev$val` es el vector de autovalores de `Sm` y `ev$vec` es la matriz de los correspondientes autovectores. Si sólo quisiéramos almacenar los autovalores podríamos haber hecho la asignación:

```
> evals ← eigen(Sm)$values
```

y en este caso `evals` sólo contendría los autovalores, habiéndose descartado la segunda componente de la lista. Si se utiliza la expresión

```
> eigen(Sm)
```

se imprimen las dos componentes, con sus nombres, en la pantalla.

---

<sup>16</sup>Aunque R permite lo contrario, deberían ser distintos entre sí

### 6.3.2 Descomposición en valores singulares. Determinantes

La función `svd` admite como argumento una matriz cualquiera,  $M$ , y calcula su descomposición en valores singulares, que consiste en obtener tres matrices,  $U$ ,  $D$  y  $V$ , tales que la primera es una matriz de columnas ortonormales con el mismo espacio de columnas que  $M$ , la segunda es una matriz diagonal de números no negativos, y la tercera es una matriz de columnas ortonormales con el mismo espacio de filas que  $M$ , tales que  $M = U \%*\% D \%*\% t(V)$ .  $D$  se devuelve en forma de vector correspondiente a los elementos diagonales. El resultado de la función es una lista de tres componentes cuyos nombres son `d`, `u` y `v`, y que corresponden a las matrices descritas.

Si  $M$  es una matriz cuadrada, es fácil ver que

```
> AbsDetM ← prod(svd(M)$d)
```

calcula el valor absoluto del determinante de  $M$ . Si necesita este cálculo a menudo, puede definirlo como una nueva función en R:

```
> AbsDet ← function(M) prod(svd(M)$d)
```

tras lo cual puede usar `AbsDet()` como cualquier otra función de R. Se deja como ejercicio, trivial pero útil, el cálculo de una función, `tr()`, que calcule la traza de una matriz cuadrada. Tenga en cuenta que no necesita realizar ninguna iteración; estudie atentamente el código de la función anterior.

### 6.3.3 Ajuste por mínimos cuadrados. Descomposición $QR$

La función `lsfit()` devuelve una lista que contiene los resultados de un ajuste por mínimos cuadrados. Una asignación de la forma

```
> MinCua ← lsfit(X, y)
```

almacena los resultados del ajuste por mínimos cuadrados de un vector de observaciones,  $y$ , y una matriz de diseño,  $X$ . Para ver más detalles puede consultar la ayuda, y también la de la función `ls.diag()` para los diagnósticos de regresión. Tenga en cuenta que no necesita incluir un término independiente en  $X$ , ya que se incluye automáticamente.

Otra función estrechamente relacionada es `qr()` y similares. Considere las siguientes asignaciones:

```
> Xplus ← qr(X)
> b ← qr.coef(Xplus, y)
> fit ← qr.fitted(Xplus, y)
> res ← qr.resid(Xplus, y)
```

que calculan la proyección ortogonal de  $y$  sobre  $X$  y la almacenan en `fit`, la proyección sobre el complemento ortogonal en `res` y el vector de coeficientes para la proyección en `b`<sup>17</sup>.

No se presupone que  $X$  sea de rango completo. Se buscan las redundancias y, si existen, se eliminan.

Esta forma es la forma antigua, a bajo nivel, de realizar ajustes de mínimos cuadrados. Aunque sigue siendo útil en algún contexto, debería ser reemplazada por los modelos estadísticos, como se verá en §10.

---

<sup>17</sup>`b` es esencialmente el resultado del operador “barra hacia atrás” de MATLAB.

## 6.4 Hojas de datos (Data frames)

Una hoja de datos<sup>18</sup> (Data frame) es una lista que pertenece a la clase `data.frame`. Hay restricciones en las listas que pueden pertenecer a esta clase, en particular:

- Los componentes deben ser vectores (numéricos, cadenas de caracteres, o lógicos), factores, matrices numéricas, listas u otras hojas de datos.
- Las matrices, listas, y hojas de datos contribuyen a la nueva hoja de datos con tantas variables como columnas, elementos o variables posean, respectivamente.
- Los vectores numéricos y los factores se incluyen sin modificar, los vectores no numéricos se fuerzan a factores cuyos niveles son los únicos valores que aparecen en el vector.
- Los vectores que constituyen la hoja de datos deben tener todos la *misma longitud*, y las matrices deben tener el mismo *tamaño de filas*.

Las hojas de datos pueden interpretarse, en muchos sentidos, como matrices cuyas columnas pueden tener diferentes modos y atributos. Pueden imprimirse en forma matricial y se pueden extraer sus filas o columnas con la indexación de matrices.

### 6.4.1 Construcción de hojas de datos

Puede construir una hoja de datos utilizando la función `data.frame`:

```
> cont <- data.frame(dom=FactorEstado, bot=ingresos, dis=FactorIngresos)
```

Puede *forzar* a que una lista cuyos componentes cumplan las restricciones para ser una hoja de datos lo sea, mediante la función `as.data.frame()`

La manera más sencilla de construir una hoja de datos es utilizar la función `read.table()` para leerla desde un archivo del sistema operativo. Esta forma se tratará en §7.

### 6.4.2 Funciones `attach()` y `detach()`

La notación `$` para componentes de listas, como por ejemplo `cont$dom`, no siempre es la más apropiada. Sería cómodo que los componentes de una lista o de una hoja de datos pudiesen ser tratados temporalmente como variables cuyo nombre fuese el del componente, sin tener que especificar explícitamente el nombre de la lista.

La función `attach()` puede tener como argumento el nombre de una lista o de una hoja de datos y permite conectar la lista o la hoja de datos directamente. Supongamos que `lentejas` es una hoja de datos con tres variables, `lentejas$u`, `lentejas$v` y `lentejas$w`. La orden

```
> attach(lentejas)
```

conecta la hoja de datos colocándola en la segunda posición de la lista de búsqueda y, supuesto que no existen variables denominadas `u`, `v` o `w` en la primera posición, `u`, `v` y `w` aparecerán como variables por sí mismas. Si realiza una asignación como por ejemplo

```
> u <- v+w
```

no se sustituye la componente `u` de la hoja de datos, sino que se crea una nueva variable `u` en el directorio de trabajo, en la primera posición de la lista de búsqueda, que enmascarará a la variable `u` de la hoja de datos. Para realizar un cambio en la propia hoja de datos, basta con utilizar la notación `$`:

```
> lentejas$u <- v+w
```

---

<sup>18</sup>Hemos utilizado esta traducción por analogía con la “hoja de cálculo”

Este nuevo valor de la componente `u` no será visible de modo directo hasta que desconecte y vuelva a conectar la hoja de datos.

Para desconectar una hoja de datos, utilice la función

```
> detach()
```

Esta función desconecta la entidad que se encuentre en la segunda posición de la lista de búsqueda. Una vez realizada esta operación dejarán de existir las variables `u`, `v` y `w` como tales, aunque seguirán existiendo como componentes de la hoja de datos.

**Nota:** La lista de búsqueda puede almacenar un número finito de elementos, por tanto no debe conectar una hoja de datos más de una vez. Del mismo modo, es conveniente desconectar la hoja de datos cuando termine de utilizar sus variables.

**Nota:** En la versión actual de R sólo se pueden conectar listas y hojas de datos en la posición 2 o superior. No es posible asignar directamente en una lista u hoja de datos conectada (por tanto, en cierto sentido, son estáticas). Es posible que cambie esta característica en el futuro.

### 6.4.3 Trabajo con hojas de datos

Una metodología de trabajo para tratar diferentes problemas utilizando el mismo directorio de trabajo es la siguiente:

- Reúna todas las variables de un mismo problema en una hoja de datos y déle un nombre apropiado e informativo;
- Para analizar un problema conecte, mediante `attach()`, la hoja de datos correspondiente (en la posición 2) y utilice el directorio de trabajo (en la posición 1) para los cálculos y variables temporales;
- Antes de terminar un análisis, añada las variables que deba conservar a la hoja de datos utilizando la forma `$` para la asignación y desconecte la hoja de datos mediante `detach()`;
- Para finalizar, elimine del directorio de trabajo las variables que no desee conservar para mantenerlo lo más limpio posible.

De este modo podrá analizar diferentes problemas utilizando el mismo directorio, aunque todos ellos compartan variables denominadas `x`, `y` o `z`, por ejemplo.

## 7 Lectura de datos de un archivo

Los datos suelen leerse desde archivos externos y no teclearse de modo interactivo. Las capacidades de lectura de archivos de R son sencillas y sus requisitos son bastante estrictos cuando no inflexibles. Se presupone que el usuario es capaz de modificar los archivos de datos con otras herramientas, por ejemplo con editores de texto<sup>19</sup>, para ajustarlos a las necesidades de R. Generalmente esta tarea es muy sencilla.

La función<sup>20</sup> `make.fields()` puede utilizarse para convertir un archivo con campos de anchura fija no delimitados en un archivo de campos delimitados. La función `count.fields()` cuenta el número de campos por línea de un archivo de campos delimitados. Estas dos funciones pueden resolver algunos problemas elementales, pero en la mayoría de los casos es mejor preparar el archivo a las necesidades de R antes de comenzar el análisis.

Si los datos se van a almacenar en hojas de datos, que es lo que recomendamos, puede leer los datos correspondientes a las mismas con la función `read.table()`. Existe también una función más genérica, `scan()`, que puede utilizar directamente.

### 7.1 La función `read.table()`

Para poder leer una hoja de datos directamente, el archivo externo debe reunir las condiciones adecuadas. La forma más sencilla es:

- La primera línea del archivo debe contener el *nombre* de cada variable de la hoja de datos.
- En cada una de las siguientes líneas, el primer elemento es la *etiqueta de la fila*, y a continuación deben aparecer los valores de cada variable.

Si el archivo tiene un elemento menos en la primera línea que en las restantes, obligatoriamente será el diseño anterior el que se utilice. En la figura 2 aparece un ejemplo de un archivo, `datos.casas`, con datos de viviendas, preparado para su lectura con esta función.

---

	Precio	Superficie	Área	Habitaciones	Años	Calefac.central
01	52.00	111.0	830	5	6.2	no
02	54.75	128.0	710	5	7.5	no
03	57.50	101.0	1000	5	4.2	no
04	57.50	131.0	690	6	8.8	no
05	59.75	93.0	900	5	1.9	si
...						

---

Figura 2: Archivo de entrada con nombres y etiquetas de filas

Predeterminadamente, los elementos numéricos (excepto las etiquetas de las filas) se almacenan como variables numéricas; y los no numéricos, como `Calefac.central`, se fuerzan como factores. Es posible modificar esta acción.

La función `read.table()` permite leer la hoja de datos directamente

```
> PreciosCasas <- read.table("datos.casas")
```

A menudo no se dispone de etiquetas de filas. En ese caso, también es posible la lectura y el programa añadirá unas etiquetas predeterminadas. Así, si el archivo tiene la forma de la figura 3, podrá leerse utilizando un parámetro adicional

<sup>19</sup>En UNIX puede utilizar el programa `perl` o los editores `sed` y `awk`. Existen versiones para MICROSOFT WINDOWS.

<sup>20</sup>Aún no está implementada esta función.



---

Precio	Superficie	Área	Habitaciones	Años	Calefac.central
52.00	111.0	830	5	6.2	no
54.75	128.0	710	5	7.5	no
57.50	101.0	1000	5	4.2	no
57.50	131.0	690	6	8.8	no
59.75	93.0	900	5	1.9	si

...

---

Figura 3: Archivo sin etiquetas de filas

```
> PreciosCasas ← read.table("datos.casas", header=T)
```

donde el parámetro adicional `header=T` indica que la primera línea es una línea de cabeceras y, por tanto, que no existen etiquetas de archivos implícitas.

## 7.2 La función `scan()`

Supongamos que el archivo `entrada.txt` contiene los datos correspondientes a tres vectores, el primero de tipo carácter y los otros dos de tipo numérico, escritos de tal modo que en cada línea aparecen los valores correspondientes de cada uno de ellos.

En primer lugar, es posible utilizar la función `scan()` para leer los tres vectores, del siguiente modo

```
> entrada ← scan("entrada.txt", list(,0,0))
```

El segundo argumento es una estructura de lista ficticia que establece el modo de los tres vectores que se van a leer. El resultado se almacena en `entrada`, y es una lista con tres componentes correspondientes a los vectores leídos. Puede referir cada uno de los vectores mediante la indexación:

```
> etiqueta ← entrada[[1]]; x ← entrada[[2]]; y ← entrada[[3]]
```

También podría haber utilizado nombres en la lista, por ejemplo

```
> entrada ← scan("entrada.txt", list(etiqueta=,x=0,y=0))
```

Puede ahora referir cada uno de los vectores con la notación `$`:

```
> etiqueta ← entrada$etiqueta; x ← entrada$x; y ← entrada$y
```

o puede conectar la lista con la función `attach` (vea §6.4.2).

Si el segundo argumento hubiese sido un sólo elemento y no una lista, todos los elementos del archivo deberían ser del tipo indicado y se hubiesen leído en un sólo vector.

```
> X ← matrix(scan("luz.txt", 0), ncol=5, byrow=T)
```

La función `scan` permite realizar lecturas más complejas, como puede consultar en la ayuda.

## 8 Ciclos. Ejecución condicional

### 8.1 Expresiones agrupadas

R es un lenguaje de expresiones, en el sentido de que el único tipo de orden que posee es una función o expresión que devuelve un resultado. Incluso una asignación es una expresión, cuyo resultado es el valor asignado<sup>21</sup> y puede utilizarse en cualquier sitio en que pueda utilizarse una expresión. En particular es posible realizar asignaciones múltiples.

Los órdenes pueden agruparse entre llaves,  $\{\text{expr}_1; \text{expr}_2; \dots; \text{expr}_m\}$ , en cuyo caso el valor del grupo es el resultado de la última expresión del grupo que se haya evaluado. Puesto que un grupo es por sí mismo una expresión, puede incluirse entre paréntesis y ser utilizado como parte de una expresión mayor. Este proceso puede repetirse si se considera necesario.

### 8.2 Órdenes de control

#### 8.2.1 Ejecución condicional: la función `if`

Existe una construcción condicional de la forma

```
> if (expr1) expr2 else expr3
```

donde  $\text{expr}_1$  debe producir un valor lógico, si éste es verdadero (T) se ejecutará  $\text{expr}_1$ . Si es falso (F) y se ha escrito la opción `else`, que es opcional, se ejecutará  $\text{expr}_3$ .

#### 8.2.2 Ciclos: Funciones `for`, `repeat` y `while`

Existe una construcción repetitiva de la forma

```
> for (nombre in expr1) expr2
```

donde *nombre* es la variable de control de iteración,  $\text{expr}_1$  es un vector (a menudo de la forma  $(n_1 : n_2)$ ), y  $\text{expr}_2$  es a menudo una expresión agrupada en cuyas sub-expresiones aparece a menudo la variable de control, *nombre*.  $\text{expr}_2$  se evalúa repetidamente conforme *nombre* recorre los valores del vector  $\text{expr}_1$ .

Por ejemplo, suponga que `ind` es un vector de indicadores de clase y se quieren hacer gráficos de `y` sobre `x`, separados para cada clase. Una posibilidad es usar la función `coplot()`, que veremos más adelante, que produce una matriz de gráficos correspondientes a cada nivel del factor. Otra forma de hacerlo es usar la función `for`:

```
> yc ← split(y, ind); xc ← split(x, ind)
> for (i in 1:length(yc)){plot(xc[[i]], yc[[i]]);
  abline(lsfilt(xc[[i]], yc[[i]]))}
```

(La función `split()` produce una lista de vectores dividiendo un vector de acuerdo a las clases especificadas por un factor. Consulte la ayuda para obtener más detalles.)

**Nota:** En R la función `for()` es más lenta que otras opciones (aunque no tanto como en S) por lo que debe evitarse cuando sea posible. Muchas funciones, como `apply()`, `tapply()`, `sapply()` y otras, están diseñadas primariamente para evitar el uso de `for()`.

Otras estructuras de repetición son

```
> repeat expresión
```

---

<sup>21</sup>La asignación devuelve el resultado de modo invisible. Basta escribirla entre paréntesis para comprobar lo dicho.

y

```
> while (condición) expresión
```

La función **break** se utiliza para terminar cualquier ciclo. Esta es la única forma (salvo que se produzca un error) de finalizar un ciclo **repeat**.

La función **next** deja de ejecutar el resto de un ciclo y pasa a ejecutar el siguiente.

Los órdenes de control se utilizan habitualmente en la escritura de *funciones*, que se tratarán en §9, donde se verán varios ejemplos.

## 9 Escritura de nuevas funciones

Como hemos visto informalmente hasta ahora, R permite crear objetos del modo *function*, que constituyen nuevas funciones de R, que se pueden utilizar a su vez en expresiones posteriores. En este proceso, el lenguaje gana considerablemente en potencia, comodidad y elegancia, y aprender a escribir funciones útiles es una de las mejores formas de conseguir que el uso de R sea cómodo y productivo.

Debemos recalcar que muchas de las funciones que se suministran con R, como `mean`, `var` o `postscript`, están de hecho escritas en R y, por tanto, no difieren materialmente de las funciones que pueda escribir el usuario.

Para definir una función debe realizar una asignación de la forma

```
> NombreDeFuncion ← function(arg1, arg2, ...) expresión
```

donde *expresión* es una expresión de R (normalmente una expresión agrupada) que utiliza los argumentos  $arg_i$  para calcular un valor que es devuelto por la función.

El uso de la función es normalmente de la forma

```
NombreDeFuncion(expr1, expr2, ...)
```

y puede realizarse en cualquier lugar en que el uso de una función sea correcto.

### 9.1 Ejemplos elementales

En primer lugar, consideremos una función que calcule el estadístico  $t$  de Student para dos muestras realizando “todos los pasos”. Este es un ejemplo muy artificial, naturalmente, ya que hay otros modos, mucho más sencillos, de obtener el mismo resultado.

La función se define del siguiente modo:

```
> DosMuestras ← function(y1, y2) {
  n1 ← length(y1); n2 ← length(y2)
  yb1 ← mean(y1); yb2 ← mean(y2)
  s1 ← var(y1); s2 ← var(y2)
  s ← ((n1-1)*s1 + (n2-1)*s2)/(n1+n2-2)
  test ← (yb1 - yb2)/sqrt(s2*(1/n1 + 1/n2))
  test
}
```

Una vez definida esta función, puede utilizarse para realizar un contraste de  $t$  de Student para dos muestras, del siguiente modo:

```
> tStudent ← DosMuestras(datos$hombre, datos$mujer); tStudent
```

En segundo lugar, considere por ejemplo la creación de una función<sup>22</sup> que calcule los coeficientes de la proyección ortogonal del vector  $y$  sobre el espacio de las columnas de la matriz  $X$ , esto es, la estimación de los coeficientes de regresión por mínimos cuadrados. Esta tarea se realizaría normalmente con la función `qr()`; sin embargo es algo compleja y compensa tener una función como la siguiente para usarla directamente.

Dado un vector,  $y^{n \times 1}$ , y una matriz,  $X^{n \times p}$ , entonces

$$X \backslash y \stackrel{\text{def}}{=} (X'X)^- X'y$$

donde  $(X'X)^-$  es la inversa generalizada de  $X'X$ . Podemos definir la función `Proyeccion` del siguiente modo

---

<sup>22</sup>En MATLAB, sería la orden `\`

```
> Proyeccion ← function(X, y) {
  X ← qr(X)
  qr.coef(X, y)
}
```

Una vez creada, puede utilizarla en cualquier expresión, como en la siguiente:

```
> CoefReg ← Proyeccion(matrizX, variabley)
```

La función `lsfit()` realiza la misma acción<sup>23</sup>. También utiliza las funciones `qr()` y `qr.coef()` en la forma anterior para realizar esta parte del cálculo. Por lo tanto puede ser interesante haber aislado esta parte en una función si se va a utilizar frecuentemente. Si ello es así, probablemente sería conveniente construir un operador binario matricial para que el uso sea más cómodo.

## 9.2 Cómo definir un operador binario

Si hubiésemos dado a la función `Proyeccion` otro nombre, de la forma

`%nombre%`

podría utilizarse como un *operador binario* en vez de con la forma funcional. Suponga, por ejemplo, que elige<sup>24</sup> ! como carácter interno. La definición de la función debería comenzar así:

```
> "%!%" ← function(X, y) {... }
```

donde hay que destacar la utilización de comillas. Una vez definida la función se utilizaría de la forma `X %!% y`.

Los operadores producto matricial, `%*%`, y producto exterior, `%o%`, son ejemplos de operadores binarios definidos de esta forma.

## 9.3 Argumentos con nombre. Valores predeterminados. Argumento “...”

Ya vimos en §2.3 que cuando los argumentos se dan por nombre, “*nombre=objeto*”, el orden es irrelevante. Además pueden utilizarse ambas formas simultáneamente: se puede comenzar dando argumentos por posición y después añadir argumentos por nombre.

Esto es, si la función `fun1` está definida como

```
> fun1 ← function(datos, hoja.datos, grafico, limite) {[aquí iría la definición]
}
```

las siguientes llamadas a la función son equivalentes:

```
> resultado ← fun1(d, hd, T, 20)
> resultado ← fun1(d, hd, grafico=T, limite=20)
> resultado ← fun1(datos=d, limite=20, grafico=T, hoja.datos=hd)
```

En muchos casos, puede suministrarse un valor predeterminado para algunos argumentos, en cuyo caso al ejecutar la función el argumento puede omitirse si el valor predeterminado es apropiado. Por ejemplo, si `fun1` estuviese definida como

```
> fun1 ← function(datos, hoja.datos, grafico=T, limite=20) {...23}
```

la llamada a la función

<sup>23</sup>Vea también los métodos descritos en §10

<sup>24</sup>El uso del símbolo `\` para el nombre, como en MATLAB, no es una elección conveniente, ya que presenta ciertos problemas en este contexto.

```
> resultado ← fun1(d, hd)
```

sería equivalente a las tres llamadas anteriores. Tenga en cuenta que puede modificar los valores predeterminados, como en el caso siguiente:

```
> resultado ← fun1(d, hd, limite=10)
```

Es importante destacar que los valores predeterminados pueden ser expresiones arbitrarias, que incluso involucren otros argumentos de la misma función, y no están restringidos a ser constantes como en el ejemplo anterior.

Otra necesidad frecuente es la de que una función pueda pasar los valores de sus argumentos a otra función. Por ejemplo, muchas funciones gráficas utilizan la función `par()` y funciones como `plot()` permiten al usuario pasar los parámetros gráficos a `par()` para controlar el resultado gráfico. (Vea §11.4.1 para detalles adicionales sobre la función `par()`.) Esta acción puede realizarse incluyendo un argumento adicional, "...", en la función, que puede ser traspasado. Un bosquejo de ejemplo se ha incluido en la figura 4.

---

```
fun1 ← function(datos, hoja.datos, grafico=T, limite=20, ...) {
  [Algunas órdenes]
  if (grafico)
    par(pch="*", ...)
  [Más órdenes]
}
```

---

Figura 4: Uso del argumento "..."

---

## 9.4 Las asignaciones dentro de una función son locales. Marco.

Es fundamental tener en cuenta que *cualquier asignación ordinaria realizada dentro de una función es local y temporal y se pierde tras salir de la función*. Por tanto, la asignación  $X \leftarrow \text{qr}(X)$  no afecta al valor del argumento de la función en que se utiliza.

Para comprender completamente las reglas que gobiernan el ámbito de las asignaciones en R es necesario familiarizarse con la noción de *marco* (frame) de evaluación. Este es un tema complejo que no será tratado en estas notas.

Si desea realizar asignaciones globales y permanentes dentro de una función, deberá utilizar el operador de 'superasignación', `<<-`, o la función `assign`. Puede encontrar una explicación más detallada consultando la ayuda.

El operador `<<-` es diferente en R y en S-PLUS. Las diferencias serán tratadas en §9.6.

## 9.5 Ejemplos más complejos

### 9.5.1 Factores de eficiencia en diseño en bloques

Estudiaremos ahora un ejemplo de función más completo: el cálculo de factores de eficiencia en un diseño en bloques. (Algunos aspectos de este problema ya han sido tratados en §5.3.)

Un diseño en bloques está definido por dos factores, por ejemplo **bloques** (**b** niveles) y **variedades**, (**v** niveles). Si  $R^{v \times v}$  y  $K^{b \times b}$  son las matrices de *réplicas* y *tamaño de bloque*, y  $N^{b \times v}$  es la matriz de incidencia, entonces los factores de eficiencia se definen como los autovalores de la matriz

$$E = I_v - R^{-1/2} N' K^{-1} N R^{-1/2} = I_v - A' A$$

donde  $A = K^{-1/2} N R^{-1/2}$ . La figura 5 contiene una forma de escribir la función.

---

```
> EfiDisBlo ← function(bloques, variedades) {
  bloques ← as.factor(bloques)           # pequeña precaución
  b ← length(levels(bloques))
  variedades ← as.factor(variedades)      # pequeña precaución
  v ← length(levels(variedades))
  K ← as.vector(table(bloques))           # elimina el atributo dim
  R ← as.vector(table(variedades))        # elimina el atributo dim
  N ← table(bloques, variedades)
  A ← 1/sqrt(K) * N * rep(1/sqrt(R), rep(b, v))
  sv ← svd(A)
  list(eficiencia=1 - sv$d^2, cvbloques=sv$u, cvvariedad=sv$v)
}
```

---

Figura 5: Una función para calcular la eficiencia de un diseño en bloques

---

Numéricamente, es levemente mejor trabajar con la función descomposición SVD en vez de con la función de los autovalores.

El resultado de esta función es una lista que contiene los factores de eficiencia como primera componente, y que además incluye dos contrastes, puesto que a veces estos suministran información adicional útil.

### 9.5.2 Cómo eliminar los nombres al imprimir una variable indexada

Para imprimir grandes matrices o variables indexadas en general, a menudo es interesante hacerlo en forma compacta sin los nombres de variables. La simple eliminación del atributo `dimnames` no es suficiente, sino que la solución consiste en asignar a dicho atributo cadenas de caracteres vacías. Por ejemplo, para imprimir la matriz `X` puede escribir

```
> temp ← X
> dimnames(temp) ← list(rep(, nrow(X)), rep(, ncol(X)))
> temp; rm(temp)
```

Este resultado puede conseguirse fácilmente definiendo la función `SinNombres` que aparece en la figura 6, que da un pequeño rodeo para conseguir el mismo resultado al tiempo que ilustra el hecho de que las funciones pueden ser cortas y al mismo tiempo muy efectivas y útiles.

Una vez definida la función, para imprimir la matriz `X` en forma compacta basta con escribir

```
> SinNombres(X)
```

Esta función es de especial utilidad al imprimir grandes variables indexadas enteras en que el interés real está más en los posibles patrones antes que en los valores en sí mismos.

### 9.5.3 Integración numérica recursiva

Las funciones pueden ser recursivas e, incluso, pueden definir funciones en su interior. Advierta, sin embargo, que dichas funciones, y por supuesto las variables, no se heredan por funciones llamadas en marcos de evaluación superior como lo serían si estuviesen en la lista de búsqueda.

---

```

SinNombres ← function(a) {
#
# Elimina los nombres de dimensiones para impresión compacta.
#
  d ← list()
  l ← 0
  for(i in dim(a)) {
    d[[l ← l + 1]] ← rep("", i)
  }
  dimnames(a) ← d
  a
}

```

Figura 6: Función de impresión de variables indexadas en forma compacta

---

El ejemplo de la figura 7 muestra una forma ingenua de realizar integración numérica unidimensional. El integrando se evalúa en los extremos del intervalo y en el centro. Si el resultado de aplicar la regla del trapecio a un solo intervalo es bastante próxima al resultado de aplicarlo a los dos, entonces este último valor se considera el resultado. En caso contrario se aplica el procedimiento a cada uno de los dos intervalos. El resultado es un proceso de integración adaptativo que concentra las evaluaciones de la función en las regiones en que es menos lineal. Conlleva, sin embargo, un gran consumo de recursos, y la función solo es competitiva con otros algoritmos cuando el integrando es tanto suave como difícil de evaluar.

El ejemplo es también un pequeño rompecabezas de programación en R.

---

```

area ← function(f, a, b, eps = 1.0e-06, lim = 10)
{
  fun1 ← function(f, a, b, fa, fb, a0, eps, lim, fun)
  { ## La función 'fun1' sólo es visible dentro de 'area'
    d ← (a + b)/2
    h ← (b - a)/4
    fd ← f(d)
    a1 ← h * (fa + fd)
    a2 ← h * (fd + fb)
    if(abs(a0 - a1 - a2) < eps || lim == 0)
      return(a1 + a2)
    else {
      return(fun(f, a, d, fa, fd, a1, eps, lim - 1, fun) +
             fun(f, d, b, fd, fb, a2, eps, lim - 1, fun))
    }
  }
  fa ← f(a)
  fb ← f(b)
  a0 ← ((fa + fb) * (b - a))/2
  fun1(f, a, b, fa, fb, a0, eps, lim, fun1)
}

```

Figura 7: Una función recursiva dentro de otra función

---



## 9.6 Ámbito

Este apartado es algo más técnico que otras partes de este documento. Sin embargo, pormenoriza una de las mayores diferencias entre S-PLUS y R.

Los símbolos que tienen lugar en el cuerpo de una función se dividen en tres clases: parámetros formales, variables locales y variables libres. Los parámetros formales son los que aparecen en la lista de argumentos de la función y sus valores quedan determinados por el proceso de asignación de los argumentos de la función a los parámetros formales. Las variables locales son aquellas cuyos valores están determinados por la evaluación de expresiones en el cuerpo de las funciones. Las variables que no son parámetros formales ni variables locales se denominan variables libres. Las variables libres se transforman en variables locales si se les asigna valor. Para aclarar los conceptos, consideremos la siguiente función:

```
f ← function(x) {
  y ← 2*x
  print(x)
  print(y)
  print(z)
}
```

En esta función, `x` es un parámetro formal, `y` es una variable local y `z` es una variable libre.

En R la asignación de valor a una variable libre se realiza consultando el entorno en el que la función se ha creado. En primer lugar definamos la función `cubo`:

```
cubo ← function(n) {
  sq ← function() n*n
  n*sq()
}
```

La variable `n` de la función `sq` no es un argumento para esta función. Por tanto es una variable libre y las reglas de ámbito deben utilizarse para determinar el valor asociado con ella. Bajo un ámbito estático es un parámetro para la función `cubo` puesto que hay una asignación activa para la variable `n` en el momento en que se define la función `sq`. La diferencia de evaluación entre R y S-PLUS es que S-PLUS intenta encontrar una variable global llamada `n` en tanto que R primero intenta encontrar una variable llamada `n` en el entorno creado cuando se activó `cubo`.

```
# primera evaluación en S
S> cubo(2)
Error in sq(): Object "n" not found
Dumped
S> n ← 3
S> cubo(2)
[1] 18
# la misma función evaluada en R
R> cubo(2)
[1] 8
```

El ámbito lexicográfico puede utilizarse para conceder a las funciones un *estado cambiante*. En el siguiente ejemplo mostramos cómo puede utilizarse R para simular una cuenta bancaria. Una cuenta bancaria necesita tener un balance o total, una función para realizar depósitos, otra para retirar fondos, y una última para conocer el balance.

Conseguiremos esta capacidad creando tres funciones dentro de `anota.importe` y devolviendo una lista que los contiene. Cuando se ejecuta `anota.importe` toma un argumento numérico, `total`, y devuelve una lista que contiene las tres funciones. Puesto que estas funciones están definidas dentro de un entorno que contiene a `total`, éstas tendrán acceso a su valor.

El operador de asignación especial, `<<-`, se utiliza para cambiar el valor asociado con `total`. Este operador comprueba los entornos creados desde el actual hasta el primero hasta encontrar uno que contenga el símbolo `total` y cuando lo encuentra, sustituye su valor en dicho entorno por el valor de la derecha de la expresión. Si se alcanza el nivel superior, correspondiente al entorno global, sin encontrar dicho símbolo, entonces lo crea en él y realiza la asignación. Para muchos usos, `<<-` crea una variable global y le asigna el valor de la derecha de la expresión<sup>24</sup>. Solo cuando `<<-` ha sido utilizado en una función que ha sido devuelta como el valor de otra función ocurrirá la conducta especial que hemos descrito.

---

```
anota.importe ← function(total) {

  list(
    deposito = function(importe) {
      if(importe <= 0)
        stop("<Los depósitos deben ser positivos!\n")
      total <<- total + importe
      cat("Depositado",importe,". El total es", total, "\n\n")
    },
    retirada = function(importe) {
      if(importe > total)
        stop("<No tienes tanto dinero!\n")
      total <<- total - importe
      cat("Descontado", importe,". El total es", total, "\n\n")
    },
    balance = function() {
      cat("El total es", total, "\n\n")
    }
  )
}

Antonio ← anota.importe(100)
Roberto ← anota.importe(200)

Antonio$retirada(30)
Antonio$balance()
Antonio$balance()

Roberto$deposit(50)
Roberto$balance()
Roberto$retirada(500)
```

Figura 8: Una función que utiliza ámbito léxico.

---

## 9.7 Adaptación del entorno

El usuario de R puede adaptar el entorno de trabajo de varias formas. Existe un archivo de inicialización del sistema y cada directorio puede tener su propio archivo de inicialización especial. Por último, puede usar las funciones especiales `.First` y `.Last`.

El archivo de inicialización del sistema se denomina `Rprofile` y se encuentra en el subdirectorio `library` del directorio inicial de R. Las órdenes contenidas en este archivo se ejecutan cada vez que se comienza una sesión de R, sea cual sea el usuario. Existe un segundo archivo, personal,

---

<sup>24</sup>En cierto sentido esto emula la conducta en S-PLUS puesto que en S-PLUS este operador siempre crea o asigna a una variable global.

denominado `.Rprofile`<sup>25</sup> que puede estar en cualquier directorio. Si ejecuta R desde un directorio que contenga este archivo, se ejecutarán las órdenes que incluya. Este archivo permite a cada usuario tener control sobre su espacio de trabajo y permite disponer de diferentes métodos de inicio para diferentes directorios de trabajo.

Si no existe el archivo `.Rprofile` en el directorio inicial, entonces R buscará si existe el archivo `.Rprofile` en el directorio inicial del usuario y, si existe, lo utilizará.

Si existe la función `.First()` en cualquiera de los dos archivos de perfil o el archivo de imagen `.RData`, recibe un tratamiento especial, ya que se ejecuta al comienzo de una sesión de R y puede utilizarse para inicializar el entorno. Por ejemplo, la definición de la figura 9 sustituye el símbolo de R para que sea `$` y establece otras características que quedan establecidas en el resto de la sesión.

En resumen, la secuencia en que se ejecutan los archivos es, `Rprofile`, `.Rprofile`, `.RData` y por último la función `.First()`. Recuerde que cualquier definición en un archivo posterior enmascarará las de un archivo precedente.

---

```
> .First ← function() {
  options(prompt="$ ", continue="+\t")
  # $ será el símbolo de sistema
  options(digits=5, length=999)
  # personaliza números y resultados
  x11()
  # abre una ventana para gráficos
  par(pch = "+")
  # carácter para realización de gráficos
  source(paste(getwd(), "/R/MisOrdenes.R", sep = ""))
  # ejecuta las órdenes contenidas en MisOrdenes.R
  library(stepfun)
  # conecta la biblioteca stepfun
}
```

---

Figura 9: Un ejemplo de función `.First`

De modo análogo, si existe la función `.Last()`, se ejecutará al término de la sesión. La figura 10 muestra un ejemplo de esta función.

---

```
> .Last ← function() {
  graphics.off()
  cat(paste(system.date(), "\nAdiós\n")) # >Es la hora de comer?
}
```

---

Figura 10: Un ejemplo de función `.Last`

## 9.8 Clases. Funciones genéricas. Orientación a objetos

La clase de un objeto determina de qué modo será tratado por lo que se conoce como funciones *genéricas*. Volviendo la oración por pasiva, una función será genérica si realiza una tarea o acción sobre sus argumentos *específica de la clase de cada argumento*. Si el argumento carece del atributo clase, o lo posee de uno no contemplado específicamente por la función genérica en cuestión, se suministra una *acción predeterminada*.

<sup>25</sup>Por tanto, en UNIX, será un archivo oculto.

El mecanismo de clase ofrece al usuario la posibilidad de diseñar y escribir funciones genéricas para propósitos especiales. Entre otras funciones genéricas puede encontrar `plot()` para representar objetos gráficamente, `summary()` para realizar análisis descriptivos de varios tipos, o `anova()` para comparar modelos estadísticos.

El número de funciones genéricas que pueden tratar una clase de modo específico puede ser muy grande. Por ejemplo, entre las funciones que pueden tratar de modo específico objetos de la clase `data.frame` se encuentran

```
[,      [[<-,   any,    as.matrix,  
[<-,   model,  plot,   summary,
```

Puede obtener una lista completa utilizando la función `methods`:

```
> methods(class="data.frame")
```

Como es esperable, el número de clases que una función genérica puede tratar también puede ser grande. Por ejemplo, la función `plot()` tiene variantes para las siguientes clases de objetos:

```
data.frame, default, density, factor,
```

y más. También en este caso puede obtener una lista completa utilizando la función `methods`:

```
> methods(plot)
```

Puede encontrar una discusión completa de este mecanismo en el manual de referencia.

## 10 Modelos estadísticos en R

En este apartado, suponemos al lector familiarizado con la terminología estadística, en particular con el análisis de regresión y el análisis de varianza. Posteriormente haremos algunas suposiciones más ambiciosas, particularmente el conocimiento de modelos lineales generalizados y regresión no lineal.

Los requisitos para el ajuste de modelos estadísticos están suficientemente bien definidos para hacer posible la construcción de herramientas generales de aplicación a un amplio espectro de problemas.

R contiene un conjunto de posibilidades que hace que el ajuste de modelos estadísticos sea muy simple. Sin embargo no llegan al nivel de **Genstat**, especialmente en cuanto a la forma del resultado que, de acuerdo con la política general de R, es más bien mínima.

### 10.1 Definición de modelos estadísticos. Fórmulas

El ejemplo básico de un modelo estadístico es un modelo de regresión lineal con errores independientes y homoscedásticos

$$y_i = \sum_{j=0}^p \beta_j x_{ij} + e_i, \quad e_i \sim \text{NID}(0, \sigma^2), \quad i = 1, 2, \dots, n$$

En notación matricial puede escribirse

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

donde  $\mathbf{y}$  es el vector de respuesta, y  $\mathbf{X}$  es la *matriz del modelo* o *matriz de diseño*, formada por las columnas  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_p$ , las variables predictoras. Muy a menudo  $\mathbf{x}_0$  será una columna de unos y definirá el *punto de corte* o *término independiente*.

#### Ejemplos.

Antes de dar una definición formal, algunos ejemplos ayudarán a centrar las ideas.

Supongamos que  $y, x, x_0, x_1, x_2, \dots$  son variables numéricas, que  $\mathbf{X}$  es una **matriz** y que  $\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots$  son factores. Las fórmulas que aparecen a continuación en la parte izquierda, especifican los modelos estadísticos descritos en la parte de la derecha.

$y \sim x$ $y \sim 1 + x$	Ambos definen el mismo modelo de regresión lineal de $y$ sobre $x$ . El primero contiene un término independiente implícito y el segundo, explícito.
$y \sim -1 + x$ $y \sim x - 1$	Regresión lineal de $y$ sobre $x$ sin término independiente.
$\log(y) \sim x_1 + x_2$	Regresión múltiple de la variable transformada, $\log(y)$ , sobre $x_1$ y $x_2$ (con un término independiente implícito).
$y \sim \text{poly}(x, 2)$ $y \sim 1 + x + \text{I}(x^2)$	Regresión polinomial de $y$ sobre $x$ de segundo grado. La primera forma utiliza polinomios ortogonales y la segunda utiliza potencias de modo explícito.
$y \sim \mathbf{X} + \text{poly}(x, 2)$	Regresión múltiple de $y$ con un modelo matricial consistente en la matriz $\mathbf{X}$ y términos polinomiales en $x$ de segundo grado.
$y \sim \mathbf{A}$	Análisis de varianza de entrada simple de $y$ , con clases determinadas por $\mathbf{A}$ .

$y \sim A + x$	Análisis de covarianza de entrada simple de $y$ , con clases determinadas por $A$ , y con covariante $x$ .
$y \sim A*B$ $y \sim A + B + A:B$ $y \sim B \%in\% A$ $y \sim A/B$	Modelo no aditivo de dos factores de $y$ sobre $A$ y $B$ . Los dos primeros especifican la misma clasificación cruzada y los dos últimos especifican la misma clasificación anidada. En términos abstractos los cuatro especifican el mismo subespacio de modelos.
$y \sim (A + B + C)^2$ $y \sim A*B*C - A:B:C$	Experimento con tres factores con un modelo que contiene efectos principales e interacciones de dos factores solamente. Ambas fórmulas especifican el mismo modelo.
$y \sim A * x$ $y \sim A/x$ $y \sim A/(1 + x) - 1$	Modelos de regresión lineal simple separados de $y$ sobre $x$ para cada nivel de $A$ . La última forma produce estimaciones explícitas de tantos términos independientes y pendientes como niveles tiene $A$ .
$y \sim A*B + \text{Error}(C)$	Un experimento con dos factores de tratamiento, $A$ y $B$ , y estratos de error determinados por el factor $C$ . Por ejemplo, un experimento split plot, con gráficos completos (y por tanto también subgráficos) determinados por el factor $C$ .

El operador  $\sim$  se utiliza para definir una *fórmula de modelo* en R. La forma, para un modelo lineal ordinario es

$$\text{respuesta} \sim [\pm] \text{term}_1 \pm \text{term}_2 \pm \text{term}_3 \pm \dots$$

**respuesta** es un vector o una matriz (o una expresión que evalúe a un vector o matriz) que definen, respectivamente, la o las variables respuesta.

$\pm$  es un operador, bien  $+$ , bien  $-$ , que implica la inclusión o exclusión, respectivamente, de un término en el modelo. El primero,  $+$ , es opcional.

**term** es un término de uno de los siguientes tipos

- una expresión vectorial, una expresión matricial, o el número 1
- un factor
- una *expresión de fórmula* consistente en factores, vectores o matrices conectados mediante *operadores de fórmula*.

En todos los casos, cada término define una colección de columnas que deben ser añadidas o eliminadas de la matriz del modelo. Un 1 significa un término independiente y predeterminadamente está incluido salvo que se elimine explícitamente.

Los *operadores de fórmula* son similares a la notación de Wilkinson y Rogers utilizada en los programas Glim y Genstat. Un cambio inevitable es que el operador “.” se ha sustituido por “:” puesto que el punto es un carácter válido para nombres de objetos en R. Un resumen de la notación se encuentra en la tabla 1 (basada en John M. Chambers and Trevor J. Hastie eds. (1992), p. 29).

Advierta que dentro de los paréntesis que habitualmente rodean los argumentos de una función todos los operadores tienen su sentido aritmético habitual. La función

$I()$  es una función identidad utilizada solamente para poder introducir términos en las fórmulas definidos utilizando operadores aritméticos.

En particular, cuando las fórmulas especifican *columnas de la matriz del modelo*, la especificación de los parámetros es implícita. Este no es el caso en otros contextos, por ejemplo en el ajuste de modelos no lineales.

Forma	Significado
$Y \sim M$	$Y$ se modeliza como $M$
$M_1 + M_2$	Incluye $M_1$ y $M_2$
$M_1 - M_2$	Incluye $M_1$ exceptuando los términos de $M_2$
$M_1 : M_2$	El producto tensorial de $M_1$ y $M_2$ . Si ambos son factores, corresponde al factor “subclases”.
$M_1 \%in\% M_2$	Similar a $M_1 : M_2$ , pero con diferente codificación.
$M_1 * M_2$	$M_1 + M_2 + M_1 : M_2$
$M_1 / M_2$	$M_1 + M_2 \%in\% M_1$
$M \sim n$	Todos los términos de $M$ junto a las “interacciones” hasta el orden $n$
$I(M)$	Aísla $M$ . Dentro de $M$ todos los operadores tienen su sentido aritmético habitual y este término aparece en la matriz del modelo.

Tabla 1: Resumen de la semántica de operadores de modelos

## 10.2 Modelos de regresión

La función primaria para el ajuste de modelos múltiples ordinarios es `lm()` y una versión resumida de su uso es la siguiente:

```
> modelo.ajustado <- lm(formula, data=hoja.de.datos)
```

Por ejemplo

```
> fm2 <- lm(y ~ x1 + x2, data=produccion)
```

ajustará un modelo de regresión múltiple de  $y$  sobre  $x_1$  y  $x_2$  (con término independiente implícito).

El término `data=produccion`, pese a ser opcional, es importante y especifica que cualquier variable necesaria para la construcción del modelo debe provenir en primer lugar de la *hoja de datos* `produccion`, y ello independientemente de que la hoja de datos `produccion` haya sido conectada a la lista de búsqueda o no.

## 10.3 Funciones genéricas de extracción de información

El valor de `lm()` es el objeto *modelo ajustado*; que consiste en una lista de resultados de clase `lm`. La información acerca del modelo ajustado puede imprimirse, extraerse, dibujarse, etc. utilizando funciones genéricas orientadas a objetos de clase `lm`. Una lista completa de las mismas es

<code>add1</code>	<code>coef</code>	<code>effects</code>	<code>kappa</code>	<code>predict</code>	<code>residuals</code>
<code>alias</code>	<code>deviance</code>	<code>family</code>	<code>labels</code>	<code>print</code>	<code>summary</code>
<code>anova</code>	<code>drop1</code>	<code>formula</code>	<code>plot</code>	<code>proj</code>	

La tabla 2 contiene una breve descripción de las más utilizadas.

## 10.4 Análisis de varianza. Comparación de modelos

El análisis de varianza, cuyo acrónimo en inglés, ANOVA, en alguna literatura en español se sustituye por ANDEVA, es otra de las técnicas aquí recogidas.

Función	Valor o efecto
<code>anova(objeto<sub>1</sub>, objeto<sub>2</sub>)</code>	Compara un submodelo con un modelo externo y produce una tabla de análisis de la varianza.
<code>coefficients(objeto)</code>	Extrae la matriz de coeficientes de regresión. Forma reducida: <code>coef(objeto)</code> .
<code>deviance(objeto)</code>	Suma de cuadrados residual, ponderada si es lo apropiado.
<code>formula(objeto)</code>	Extrae la fórmula del modelo.
<code>plot(objeto)</code>	Crea dos gráficos, uno de observado frente a predicho, el otro de residuos absolutos frente a predicho.
<code>predict(objeto, newdata=hoja.de.datos)</code>  <code>predict.gam(objeto, newdata=hoja.de.datos)</code>	La nueva hoja de datos que se indica debe tener variables cuyas etiquetas coincidan con las de la original. El resultado es un vector o matriz de valores predichos correspondiente a los valores de las variables de <i>hoja.de.datos</i> . <code>predict.gam()</code> es una forma alternativa de <code>predict()</code> que puede usar para objetos ajustados de clases <code>lm</code> , <code>glm</code> y <code>gam</code> . Esta función debe usarse, por ejemplo, en los casos en que se utilizan polinomios ortogonales como las funciones originales de ajuste y, por tanto, la adición de nuevos datos implica la utilización de funciones distintas de las originales.
<code>print(objeto)</code>	Imprime una versión concisa del objeto. A menudo se utiliza implícitamente.
<code>residuals(objeto)</code>	Extrae la matriz de residuos, ponderada si es necesario. La forma reducida es <code>resid(object)</code> .
<code>summary(objeto)</code>	Imprime un resumen estadístico completo de los resultados del análisis de regresión.

Tabla 2: Funciones genéricas utilizadas habitualmente con objetos de la clase `lm`

La función de ajuste de modelo `aov(formula, data= hoja.de.datos)` opera en el nivel más simple de modo muy similar a la función `lm()`, y muchas de las funciones genéricas contenidas en la tabla 2 le son de aplicación.

Debemos destacar que, además, `aov()` realiza un análisis de modelos estratos de error múltiple, tales como experimentos split plot, o diseños en bloques incompletos balanceados con recuperación de información inter-bloques. La fórmula

$$\text{response} \sim \text{mean.formula} + \text{Error(formula.estratos)}$$

especifica un experimento multiestrato con estratos de error definidos por *formula.estratos*. En el caso más sencillo, *formula.estratos* es un factor, cuando define un experimento con dos estratos, esto es, dentro de y entre los niveles de un factor.

Por ejemplo, con todas las variables predictoras factores, un modelo como el siguiente:

```
> mf <- aov(cosecha ~ v + n*p*k + Error(granjas/bloques), data=datos.granjas)
```

sería utilizable para describir un experimento con media del modelo `v + n*p*k` y tres estratos de error: “entre granjas”, “dentro de granjas, entre bloques” y “dentro de bloques”.

#### 10.4.1 Tablas ANOVA

Debe tenerse en cuenta que la tabla ANOVA corresponde a una sucesión de modelos ajustados. Las sumas de cuadrados que en ella aparecen corresponden a la disminución en las sumas de



cuadrados residuales resultantes de la inclusión de *un término* concreto en *un lugar* concreto de la sucesión. Por tanto el orden de inclusión sólo será irrelevante en experimentos ortogonales.

Para experimentos multiestrato el procedimiento consiste en primer lugar en proyectar la respuesta sobre los estratos de error, una vez más en secuencia, y ajustar la media del modelo a cada proyección. Para más detalles, consulte John M. Chambers and Trevor J. Hastie eds. (1992), §5.

Una alternativa más flexible a la tabla ANOVA completa es comparar dos o más modelos directamente utilizando la función `anova()`.

```
> anova(fitted.model.1, fitted.model.2, ...)
```

El resultado es una tabla ANOVA que muestra las diferencias entre los modelos ajustados cuando se ajustan en ese orden preciso. Los modelos ajustados objeto de comparación constituyen por tanto una sucesión jerárquica. Este resultado no suministra información distinta a la del caso predeterminado, pero facilita su comprensión y control.

## 10.5 Actualización de modelos ajustados. Uso de “.”

La función `update()` se utiliza muy a menudo para ajustar un modelo que difiere de uno ajustado previamente en unos pocos términos que se añaden o se eliminan. Su forma es:

```
> modelo.nuevo ← update(modelo.anterior, fórmula.nueva)
```

En *fórmula.nueva* puede utilizarse un punto, “.”, para hacer referencia a “la parte correspondiente de la fórmula del modelo anterior”. Por ejemplo,

```
> mf05 ← lm(y ~ x1 + x2 + x3 + x4 + x5, data=produccion)
> mf6 ← update(mf05, . ~ . + x6)
> rmf6 ← update(mf6, sqrt(.) ~ .)
```

ajusta una regresión múltiple con cinco variable (posiblemente) de la hoja de datos `produccion`, ajusta un modelo adicional incluyendo un sexto regresor, y ajusta una variante del modelo donde se aplica una transformación raíz cuadrada a la variable predicha.

Advierta especialmente que si se especifica el argumento `data=` en la llamada original a la función de ajuste del modelo, esta información se pasa a su vez a través del objeto modelo ajustado a la función `update()` y sus asociados.

El nombre “.” puede utilizarse también en otros contextos, pero con un significado levemente distinto. Por ejemplo,

```
> gra.completo ← lm(y ~ . , data=produccion)
```

ajustará un modelo con respuesta `y` y con variables predictoras, *todas las de la hoja de datos `produccion`*.

Otras funciones que permiten explorar sucesiones crecientes de modelos son `add1()`, `drop1()`, `step()` y `stepwise()`. Consulte la ayuda para obtener información de las mismas.

## 10.6 Modelos lineales generalizados. Familias

Los modelos lineales generalizados constituyen una extensión de los modelos lineales para tomar en consideración tanto distribuciones de respuestas no normales como transformaciones para conseguir linealidad, de una forma directa. Un modelo lineal generalizado puede describirse según las siguientes suposiciones:

- Existe una variable respuesta,  $y$ , de interés y una variables estímulo,  $x_1, x_2, \dots$ ; cuyos valores influyen en la distribución de la respuesta.

- Las variables estímulo influyen en la distribución de  $y$  mediante una *función lineal solamente*. Esta función lineal recibe el nombre de *predictor lineal*, y se escribe habitualmente

$$\eta = \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p$$

por tanto  $x_i$  no influye en la distribución de  $y$  si y sólo si  $\beta_i = 0$ .

- La distribución de  $y$  es de la forma

$$f_Y(y; \mu, \varphi) = \exp \left[ \frac{A}{\varphi} \{y\lambda(\mu) - \gamma(\lambda(\mu))\} + \tau(y, \varphi) \right]$$

donde  $\varphi$  es un *parámetro de escala* (posiblemente conocido) que permanece constante para todas las observaciones,  $A$  representa una ponderación a priori que se supone conocida pero que puede variar con las observaciones, y  $\mu$  es la media de  $y$ . Por tanto, se supone que la distribución de  $y$  queda determinada por su media y tal vez un parámetro de escala.

- La media,  $\mu$ , es una función inversible del predictor lineal:

$$\mu = m(\eta), \quad \eta = m^{-1}(\mu) = \ell(\mu)$$

y esta función inversa,  $\ell(\cdot)$  se denomina *función de enlace*.

Estas suposiciones son suficientemente amplias para acomodar una amplia clase de modelos útiles en la práctica estadística y, al tiempo, suficientemente estrictas como para permitir el desarrollo de una metodología unificada de estimación e inferencia, al menos aproximadamente. Los interesados en este tema pueden consultar cualquiera de los trabajos de referencia sobre este tema, tales como Peter McCullagh and John A. Nelder (1989()) y Annette J. Dobson (1990).

### 10.6.1 Familias

La clase de modelos lineales generalizados que pueden ser tratados en R incluye las distribuciones de respuesta *gaussian* (normal), *binomial*, *poisson*, *inverse gaussian* (normal inversa) y *gamma* así como los modelos de *quasi-likelihood* (cuasi-verosimilitud) cuya distribución de respuesta no está explícitamente definida. En este último caso debe especificarse la *función de varianza* como una función de la media, pero en el resto de casos esta función está implícita en la distribución de respuesta.

Cada distribución de respuesta admite una variedad de funciones de enlace para conectar la media con el predictor lineal. La tabla 3 recoge las que están disponibles automáticamente.

Función de enlace	Nombre de la familia					
	binomial	gaussian	Gamma	inverse.gaussian	poisson	quasi
logit	*					*
probit	*					*
cloglog	*					*
identity		*	*		*	*
inverse			*			*
log			*		*	*
1/mu^2				*		*
sqrt					*	*

Tabla 3: Familias y funciones de enlace que pueden utilizar

La combinación de una distribución de respuesta, una función de enlace y otras informaciones que son necesarias para llevar a cabo la modelización se denomina *familia* del modelo lineal generalizado.

### 10.6.2 La función `glm`

Puesto que la distribución de respuesta depende de las variables de estímulo *solamente* a través de una función lineal, se puede utilizar el mismo mecanismo de los modelos lineales para especificar la parte lineal de un modelo generalizado. Sin embargo la familia debe especificarse de modo distinto.

La función `glm()` permite ajustar un modelo lineal generalizado y tiene la forma siguiente

```
> modelo.ajustado <- glm(formula, family = generador.familia, data=hoja.datos)
```

La única característica nueva es *generador.de.familia* que es el instrumento mediante el que se describe la familia. Es el nombre de una función que genera una lista de funciones y expresiones que juntas definen y controlan el modelo y el proceso de estimación. Aunque puede parecer complicado a primera vista, su uso es bastante sencillo.

Los nombres de los generadores de familias estándar suministrados con R aparecen en la tabla 3 con la denominación de “Nombre de la familia”. Si además debe seleccionar una función de enlace, debe indicarla como un parámetro, entre paréntesis, del nombre de la familia. En el caso de la familia *quasi*, la función de varianza puede especificarse del mismo modo.

Veamos algunos ejemplos.

#### La familia gaussiana

Una expresión de la forma

```
> mf <- glm(y ~ x1+x2, family=gaussian, data=ventas)
```

obtiene el mismo resultado que

```
> mf <- lm(y ~ x1+x2, data=ventas)
```

pero con menor eficiencia. Tenga en cuenta que la familia gaussiana no dispone automáticamente de una serie de funciones de enlace, por lo que no admite parámetros. Si en un problema necesita utilizar la familia gaussiana con un enlace no estándar, la solución pasa por el uso de la familia *quasi*, como veremos posteriormente.

#### La familia binomial

Consideremos el siguiente ejemplo artificial.

Los hombres de la isla griega de Kalythos sufren una enfermedad ocular congénita cuyos efectos se acrecientan con la edad. Se tomó una muestra de varios isleños de diferentes edades cuyos resultados se muestran en la tabla 4.

Edad:	20	35	45	55	70
No. sujetos	50	50	50	50	50
No. invidentes	6	17	26	37	44

Tabla 4: Datos de ceguera en Kalythos

Consideramos el problema de ajustar un modelo logístico y otro probit a estos datos, y estimar en cada modelo el parámetro LD50, correspondiente a la edad en que la probabilidad de ceguera es del 50%.

Si  $y$  es el número de invidentes a la edad  $x$  y  $n$  es el número de sujetos estudiados, ambos modelos tienen la forma

$$y \sim B(n, F(\beta_0 + \beta_1 x))$$

donde, para el caso probit,  $F(z) = \Phi(z)$  es la función de distribución normal (0,1), y en el caso logit (que es el predeterminado),  $F(z) = e^z / (1 + e^z)$ . En ambos caso, LD50 se define como

$$LD50 = -\beta_0 / \beta_1$$

Esto es, el punto en que el argumento de la función de distribución es cero.

El primer paso es preparar los datos en una hoja de datos

```
> kalythos <- data.frame(x=c(20,35,45,55,70), n=rep(50,5),
  y=c(6,17,26,37,44))
```

Para ajustar un modelo binomial utilizando `glm()` existen dos posibilidades para la respuesta:

- Si la respuesta es un *vector*, entonces debe corresponder a datos *binarios* y por tanto sólo debe contener ceros y unos.
- Si la respuesta es una *matriz de dos columnas*, la primera columna debe contener el número de éxitos y la segunda el de fracasos.

Aquí vamos a utilizar la segunda de estas convenciones, por lo que debemos añadir una matriz a la hoja de datos:

```
> kalythos$Ymat <- cbind(kalythos$y, kalythos$n - kalythos$y)
```

Para ajustar los modelos utilizamos

```
> mfp <- glm( Ymat~x, family=binomial(link=probit), data=kalythos)
> mfl <- glm( Ymat~x, family=binomial, data=kalythos)
```

Puesto que la función de enlace logit es la predeterminada, este parámetro puede omitirse en la segunda expresión. Para ver los resultados de cada ajuste usaremos

```
> summary(mfp)
> summary(mfl)
```

Ambos modelos se ajustan bien. Para estimar LD50 podemos usar la siguiente función:

```
> ld50 <- function(b) -b[1]/b[2]
> ldp <- ld50(coef(mfp)); ldl <- ld50(coef(mfl)); c(ldp, ldl)
```

y obtendremos los valores 43.663 y 43.601 respectivamente.

## Modelos de Poisson

Para la familia poisson el enlace predeterminado es `log`, y en la práctica el uso fundamental de esta familia es ajustar modelos loglineales de Poisson a datos de frecuencias, cuya distribución en sí es a menudo multinomial. Este es un tema amplio e importante que no discutiremos aquí y que incluso constituye una parte fundamental de la utilización de modelos generalizados no gaussianos.

A veces surgen datos auténticamente poissonianos que, en el pasado se analizaban a menudo como datos gaussianos tras aplicarles una transformación logarítmica o de raíz cuadrada. Como alternativa al último, puede ajustarse un modelo lineal generalizado de Poisson, como en el siguiente ejemplo:

```
> fmod <- glm(y ~ A+B + x, family=poisson(link=sqrt), data=frec.gusanos)
```

### Modelos de cuasi-verosimilitud

En todas las familias, la varianza de la respuesta dependerá de la media, y el parámetro de escala actuará como un multiplicador. La forma de dependencia de la varianza respecto de la media es una característica de la distribución de respuesta, por ejemplo para la distribución poisson será  $\text{Var}[y] = \mu$ .

Para estimación e inferencia de cuasi-verosimilitud, la distribución de respuesta precisa no está especificada, sino más bien sólo la función de enlace y la forma en que la función de varianza depende de la media. Puesto que la estimación de cuasi-verosimilitud utiliza formalmente las mismas técnicas de la distribución gaussiana, esta familia permite ajustar modelos gaussianos con funciones de enlace no estándar o incluso con funciones de varianza.

Por ejemplo, consideremos la regresión no lineal siguiente

$$y = \frac{\theta_1 z_1}{z_2 - \theta_2} + e \quad (1)$$

que puede escribirse también de la forma

$$y = \frac{1}{\beta_1 x_1 + \beta_2 x_2} + e$$

donde  $x_1 = z_2/z_1$ ,  $x_2 = -1/x_1$ ,  $\beta_1 = 1/\theta_1$  y  $\beta_2 = \theta_2/\theta_1$ . Suponiendo que existe una hoja de datos, *bioquimica*, apropiada podemos ajustar este modelo mediante

```
> ajuste.nl <- glm(y~x1+x2-1,family=
  quasi(link=inverse,variance=constant), data=bioquimica)
```

Si desea mayor información, lea las ayudas correspondientes.

## 10.7 Modelos de Mínimos cuadrados no lineales y de Máxima verosimilitud

Ciertas formas de modelos no lineales pueden ajustarse mediante modelos lineales generalizados, con la función `glm()`, pero en la mayoría de los casos será necesario utilizar optimización no lineal. La función que lo realiza en R es `nlm()`, que reemplaza aquí a la función `ms()` de S-PLUS. Buscamos los valores de los parámetros que minimizan algún índice de falta de ajuste y `nlm()` lo resuelve probando varios parámetros iterativamente. Al contrario que en la regresión lineal, por ejemplo, no existe garantía de que el procedimiento converja a unos estimadores satisfactorios. La función `nlm()` necesita unos valores iniciales de los parámetros a estimar y la convergencia depende críticamente de la calidad de dichos valores iniciales.

### 10.7.1 Mínimos cuadrados

Una forma de ajustar un modelo no lineal es minimizar la suma de los cuadrados de los errores o residuos (SSE). Este método tiene sentido si los errores observados pueden proceder de una distribución normal.

Presentamos un ejemplo debido a D. M. Bates and D. G. Watts (1988), p. 51. Los datos son:

```
> x <- c(0.02, 0.02, 0.06, 0.06, 0.11, 0.11, 0.22, 0.22, 0.56, 0.56, 1.10, 1.10)
> y <- c(76, 47, 97, 107, 123, 139, 159, 152, 191, 201, 207, 200)
```

El modelo a ajustar es:

```
> fn <- function(p) sum((y - (p[1] * x)/(p[2] + x))^2)
```

Para realizar el ajuste necesitamos unos valores iniciales de los parámetros. Una forma de encontrar unos valores iniciales apropiados es representar gráficamente los datos, conjeturar unos valores de los parámetros, y dibujar sobre los datos la curva correspondiente a estos valores.

```
> plot(x,y)
> xajustado <- seq(.02,1.1,.05)
> yajustado <- 200*xajustado/(.1+xajustado)
> lines(spline(xajustado,yajustado))
```

Aunque se podría tratar de encontrar unos valores mejores, los valores obtenidos, 200 y .1, parecen adecuados. Ya podemos realizar el ajuste:

```
> resultado <- nlm(fn,p=c(200,.1),hessian=T)
```

Tras el ajuste, `resultado$minimum` es SSE, y `resultado$estimates` son los estimadores por mínimos cuadrados de los parámetros. Para obtener los errores típicos aproximados de los estimadores (SE) hacemos lo siguiente:

```
> sqrt(diag(2*resultado$minimum/(length(y)-2)*solve(resultado$hessian)))
```

El número 2 en dicha expresión representa el número de parámetros. Un intervalo de confianza al 95% será el estimador del parámetro  $\pm 1.96$  SE. Podemos representar el ajuste en un nuevo gráfico:

```
> plot(x,y)
> xajustado <- seq(.02,1.1,.05)
> yajustado <- 212.68384222*xajustado/(0.06412146+xajustado)
> lines(spline(xajustado,yajustado))
```

### 10.7.2 Máxima verosimilitud

El método de máxima verosimilitud es un método de ajuste de un modelo no lineal que puede aplicarse incluso cuando los errores no son normales. El método busca los valores de los parámetros que maximizan la log-verosimilitud o, lo que es igual, minimizan la -log-verosimilitud. El siguiente ejemplo, tomado de Annette J. Dobson (1990), pp. 108-11, ajusta un modelo logístico a los datos de dosis-respuesta, que claramente también podría ajustarse utilizando la función `glm()`. Los datos son:

```
> x <- c(1.6907,1.7242,1.7552,1.7842,1.8113,1.8369,1.8610,1.8839)
> y <- c(6,13,18,28,52,53,61,60)
> n <- c(59,60,62,56,63,59,62,60)
```

La -log-verosimilitud a minimizar es:

```
> fn <- function(p)
  sum( - (y*(p[1]+p[2]*x) - n*log(1+exp(p[1]+p[2]*x)) + log(choose(n,y)) ))
```

Elegimos unos valores iniciales y realizamos el ajuste:

```
> resultado <- nlm(fn,p=c(-50,20),hessian=T)
```

Tras lo cual, `resultado$minimum` es la -log-verosimilitud, y `resultado$estimates` son los estimadores de máxima verosimilitud. Para obtener los SE aproximados de los parámetros, escribimos:

```
> sqrt(diag(solve(resultado$hessian)))
```

El intervalo de confianza al 95% es, el estimador del parámetro  $\pm 1.96$  SE.

## 10.8 Algunos modelos no-estándar

Concluimos esta parte con una breve mención a otras posibilidades de R para regresión especial y análisis de datos.

**Modelos mezclados** La biblioteca creada por usuarios **nlme** contiene las funciones **lme** y **nlme** para modelos de efectos mezclados lineales y no lineales, esto es, regresiones lineales y no lineales en las cuales algunos coeficientes corresponden a efectos aleatorios. Estas funciones hacen un uso intensivo de las fórmulas para especificar los modelos.

**Regresión con aproximación local** La función **loess()** ajusta una regresión no paramétrica utilizando regresión polinomial localmente ponderada. Este tipo de regresión es útil para poner de relieve una tendencia en datos confusos o para reducir datos y obtener alguna luz de la estructura de grandes conjuntos de datos.

**Regresión robusta** Existen varias funciones para ajustar modelos de regresión resistentes a la influencia de valores anómalos (outliers) en los datos. La función **lqs** en la biblioteca del mismo nombre contiene los algoritmos para ajustes altamente resistentes. Otras funciones menos resistentes pero más eficientes estadísticamente se encuentran en otras bibliotecas creadas por usuarios, por ejemplo la función **rlm** en la biblioteca **MASS**.

**Modelos aditivos** Esta técnica intenta construir una función de regresión a partir de funciones aditivas suaves de las variables predictoras, habitualmente una por cada variable predicha. Las funciones **avas** y **ace** en la biblioteca **acepack** y las funciones **bruto** y **mars** de la biblioteca **mda** son ejemplos de estas técnicas contenidas en bibliotecas creadas por usuarios para R.

**Modelos basados en árboles** En vez de buscar un modelo lineal global explícito para predicción o interpretación, los modelos basados en árboles intentan bifurcar los datos, recursivamente, en puntos críticos de las variables predictoras con la finalidad de conseguir una partición de los datos en grupos tan homogéneos dentro del grupo y tan heterogéneos de un grupo a otro, como sea posible. Los resultados a menudo conducen a una comprensión de los datos que otros métodos de análisis de datos no suelen suministrar.

Los modelos se especifican en la forma de un modelo lineal ordinario. La función de ajuste es **tree()**, y muchas funciones genéricas, como **plot()** y **text()** pueden mostrar los resultados de este ajuste de modo gráfico.

Puede encontrar estos modelos en las bibliotecas creadas por usuarios **rpart** y **tree**.

## 11 Procedimientos gráficos

Las posibilidades gráficas son un componente de R muy importante y versátil. Es posible utilizarlas para mostrar una amplia variedad de gráficos estadísticos y también para construir nuevos tipos de gráficos.

Los gráficos pueden usarse tanto en modo interactivo como no interactivo, pero en la mayoría de los casos el modo interactivo es más productivo. Además al iniciar R en este modo, se activa un manejador para mostrar gráficos. Aunque este paso es automático, es útil conocer que la orden es `X11()`, aunque también puede usar `windows()` en MICROSOFT WINDOWS.

Las órdenes gráficas se dividen en tres grupos básicos:

**Alto nivel** Son funciones que crean un nuevo gráfico, posiblemente con ejes, etiquetas, títulos, etc..

**Bajo nivel** Son funciones que añaden información a un gráfico existente, tales como puntos adicionales, líneas y etiquetas.

**Interactivas** Son funciones que permiten interactuar con un gráfico, añadiendo o eliminando información, utilizando un dispositivo apuntador, como un ratón.

Además, en R existe una lista de *parámetros gráficos* que pueden utilizarse para adaptar los gráficos.

### 11.1 Funciones gráficas de nivel alto

Las órdenes gráficas de nivel alto están diseñadas para generar un gráfico completo a partir de unos datos pasados a la función como argumento. Cuando es necesario se generan automáticamente ejes, etiquetas o títulos (salvo que se especifique lo contrario). Estas órdenes comienzan siempre un nuevo gráfico, borrando el actual si ello es necesario.

#### 11.1.1 La función plot

Una de las funciones gráficas más utilizadas en R es `plot`, que es una función *genérica*, esto es, el tipo de gráfico producido es dependiente de la *clase* del primer argumento.

<code>plot(x,y)</code> <code>plot(xy)</code>	Si <code>x</code> e <code>y</code> son vectores, <code>plot(x,y)</code> produce un diagrama de dispersión de <code>y</code> sobre <code>x</code> . El mismo efecto se consigue suministrando un único argumento (como se ha hecho en la segunda forma) que sea bien una lista con dos elementos, <code>x</code> e <code>y</code> , bien una matriz con dos columnas.
<code>plot(x)</code>	Si <code>x</code> es una serie temporal, produce un gráfico temporal, si <code>x</code> es un vector numérico, produce un gráfico de sus elementos sobre el índice de los mismos, y si <code>x</code> es un vector complejo, produce un gráfico de la parte imaginaria sobre la real de los elementos del vector.
<code>plot(f)</code> <code>plot(f,y)</code>	Sean <code>f</code> un factor, e <code>y</code> un vector numérico. La primera forma genera un diagrama de barras de <code>f</code> ; la segunda genera diagramas de cajas de <code>y</code> para cada nivel de <code>f</code> .
<code>plot(hd)</code> <code>plot(~ expr)</code> <code>plot(y ~ expr)</code>	Sean <code>hd</code> una hoja de datos, <code>y</code> un objeto cualquiera, y <code>expr</code> una lista de nombres de objetos separados por símbolos '+' (por ejemplo, <code>a + b + c</code> ). Las dos primeras formas producen diagramas de todas las parejas de variables de la hoja de datos <code>hd</code> (en el primer caso) y de los objetos de la expresión <code>expr</code> (en el segundo caso). La tercera forma realiza sendos gráficos de <code>y</code> sobre cada objeto de la expresión <code>expr</code> (uno para cada objeto).



### 11.1.2 Representación de datos multivariantes

R posee dos funciones muy útiles para representar datos multivariantes. Si  $X$  es una matriz numérica o una hoja de datos, la orden

```
> pairs(X)
```

produce una matriz de gráficos de dispersión para cada pareja de variables definidas por las columnas de  $X$ , esto es, cada columna de  $X$  se representa frente a cada una de las demás columnas, y los  $n(n - 1)$  gráficos se presentan en una matriz de gráficos.

Cuando se trabaja con tres o cuatro variables, la función `coplot` puede ser más apropiada. Si  $a$  y  $b$  son vectores numéricos y  $c$  es un vector numérico o un factor (todos de la misma longitud) entonces la orden

```
> coplot(a ~ b | c)
```

produce diagramas de dispersión de  $a$  sobre  $b$  para cada valor de  $c$ . Si  $c$  es un factor, esto significa que  $a$  se representa sobre  $b$  para cada nivel de  $c$ . Si  $c$  es un vector numérico, entonces se agrupa en *intervalos* y para cada intervalo se representa  $a$  sobre  $b$  para los valores de  $c$  dentro del intervalo. El número y tamaño de los intervalos puede controlarse con el argumento `given.values=` de la función `coplot()`. La función `co.intervals()` también es útil para seleccionar intervalos. Asimismo, es posible utilizar dos variables *condicionantes* con una orden como

```
> coplot(a ~ b | c + d)
```

que produce diagramas de  $a$  sobre  $b$  para cada intervalo de condicionamiento de  $c$  y  $d$ .

Las funciones `coplot()` y `pairs()` utilizan el argumento `panel=` para personalizar el tipo de gráfico que aparece en cada panel o recuadro. El valor predeterminado es `points()` para producir un diagrama de dispersión, pero si introducen otras funciones gráficas de nivel bajo de los dos vectores  $x$  e  $y$  como valor de `panel=` se produce cualquier tipo de gráfico que se desee. Una función útil en este contexto es `panel.smooth()`.

### 11.1.3 Otras representaciones gráficas

Existen otras funciones gráficas de nivel alto que producen otros tipos de gráficos. Algunas de ellas son las siguientes:

---

<code>tsplot(x1,x2,...)</code>	Dibuja varias series temporales utilizando la misma escala. El uso de una escala simultánea también es útil cuando $x_i$ son vectores numéricos ordinarios, en cuyo caso se representan sobre los números 1, 2, 3, ...
--------------------------------	--

---

<code>qqnorm(x)</code>	Gráficos de comparación de distribuciones. El primero representa el vector $x$ sobre los valores esperados normales. El segundo le añade una recta que pasa por los cuartiles de la distribución y de los datos. El tercero representa los cuantiles de $x$ sobre los de $y$ para comparar sus distribuciones respectivas.
<code>qqline(x)</code>	
<code>qqplot(x,y)</code>	

---

<code>hist(x)</code>	Produce un histograma del vector numérico $x$ . El número de clases se elige habitualmente de modo correcto, pero puede tomarse una decisión con el argumento <code>nclass=</code> , o bien especificando los puntos de corte con el argumento <code>breaks=</code> . Si está presente el argumento <code>probability=T</code> , las barras representan frecuencias relativas en vez de absolutas.
<code>hist(x,nclass=n)</code>	
<code>hist(x, breaks=...)</code>	

---

---

<code>dotplot(x,...)</code>	Construye un gráfico de puntos de <b>x</b> . En este tipo de gráficos, el eje <i>y</i> etiqueta los datos de <b>x</b> y el eje <i>x</i> da su valor. Por ejemplo, permite una selección visual sencilla de todos los elementos con valores dentro de un rango determinado.
<code>image(x, y, z, breaks=...)</code> <code>contour(x, y, z, breaks=...)</code> <code>persp(x, y, z, breaks=...)</code>	Gráficos tridimensionales. <b>image</b> representa una retícula de rectángulos con colores diferentes según el valor de <b>z</b> , <b>contour</b> representa curvas de nivel de <b>z</b> , y <b>persp</b> representa una superficie tridimensional de <b>z</b> .

---

#### 11.1.4 Argumentos de las funciones gráficas de nivel alto

Existe una serie de argumentos que pueden pasarse a las funciones gráficas de nivel alto, entre otras las siguientes:

---

<code>add=T</code>	Obliga a la función a comportarse como una función a nivel bajo, de modo que el gráfico que genere se superpondrá al gráfico actual, en vez de borrarlo (sólo en algunas funciones)
<code>axes=F</code>	Suprime la generación de ejes. Útil para crear ejes personalizados con la función <b>axis</b> . El valor predeterminado es <code>axes=T</code> , que incluye los ejes.
<code>log="x"</code> <code>log="y"</code> <code>log="xy"</code>	Hace que el eje <i>x</i> , el eje <i>y</i> , o los dos ejes, sean logarítmicos. En algunos gráficos no tiene efecto.
<code>type=</code> <code>type="p"</code> <code>type="l"</code> <code>type="b"</code> <code>type="o"</code> <code>type="h"</code> <code>type="s"</code> <code>type="S"</code> <code>type="n"</code>	Este argumento controla el tipo de gráfico producido Dibuja puntos individuales. Este es el valor predeterminado Dibuja líneas Dibuja ambos: puntos y líneas que los unen Dibuja puntos y líneas que los unen cubriéndolos Dibuja líneas verticales desde cada punto al eje X Dibuja un gráfico de escalera. En la primera forma, la escalera comienza hacia la derecha, en la segunda, hacia arriba. No se realiza ningún gráfico, aunque se dibujan los ejes (salvo indicación en contra) y se prepara el sistema de coordenadas de acuerdo a los datos. Suele utilizarse para crear gráficos en los que a continuación se utilizarán órdenes de nivel bajo.
<code>xlab="cadena"</code> <code>ylab="cadena"</code>	Definen las etiquetas de los ejes <i>x</i> e <i>y</i> , en vez de utilizar las etiquetas predeterminadas, que normalmente son los nombres de los objetos utilizados en la llamada a la función gráfica de nivel alto.
<code>main="cadena"</code> <code>sub="cadena"</code>	Título del gráfico, aparece en la parte superior con tamaño de letra grande. Subtítulo del gráfico, aparece debajo del eje <i>x</i> en un tamaño de letra menor.

---

## 11.2 Funciones gráficas de nivel bajo

A veces las funciones gráficas de nivel alto no producen exactamente el tipo de gráfico deseado. En este caso pueden añadirse funciones gráficas de nivel bajo para añadir información adicional (tal como puntos, líneas o texto) al gráfico actual.

Algunas de las funciones gráficas de nivel bajo más usuales son:

<code>points(x,y)</code> <code>lines(x,y)</code>	Añade puntos o líneas conectadas al gráfico actual. El argumento <code>type=</code> de la función <code>plot()</code> puede pasarse a estas funciones (y su valor predeterminado es "p" para <code>points</code> y "l" para <code>lines</code> .)
<code>text(x, y, etiquetas, ...)</code>	<p>Añade texto al gráfico en las coordenadas <code>x</code>, <code>y</code>. Normalmente, <code>etiquetas</code> es un vector de enteros o de caracteres, en cuyo caso, <code>etiquetas[i]</code> se dibuja en el punto <code>(x[i], y[i])</code>. El valor predeterminado es <code>1:length(x)</code>.</p> <p>Nota: Esta función se utiliza a menudo en la secuencia</p> <pre>&gt; plot(x, y, type="n"); text(x, y,nombres)</pre> <p>El parámetro gráfico <code>type="n"</code> suprime los puntos pero construye los ejes, y la función <code>text</code> permite incluir caracteres especiales para representar los puntos, como se especifica en el vector de caracteres <code>nombres</code>.</p>
<code>abline(a, b)</code> <code>abline(h=y)</code> <code>abline(v=x)</code> <code>abline(lm.obj)</code>	Añade al gráfico actual, una recta de pendiente <code>b</code> y ordenada en el origen <code>a</code> . La forma <code>h=y</code> representa una recta horizontal de altura <code>y</code> , y la forma <code>v=x</code> , una similar, vertical. En el cuarto caso, <code>lm.obj</code> puede ser una lista con un componente <code>\$coefficients</code> de longitud 2 (como el resultado de una función de ajuste de un modelo) que se interpretan como ordenada y pendiente, en ese orden.
<code>polygon(x, y, ...)</code>	Añade al gráfico actual un polígono cuyos vértices son los elementos de <code>(x,y)</code> ; (opcionalmente) sombreado con líneas, o relleno de color si el periférico lo admite.
<code>legend(x,y, letrero,...)</code>	Añade al gráfico actual un letrero en la posición especificada. Los caracteres para dibujar, los estilos de líneas, los colores, etc. están identificados con los elementos del vector <code>letrero</code> . Debe darse al menos otro argumento más, <code>v</code> , un vector de la misma longitud que <code>letrero</code> , con los correspondientes valores de dibujo, como sigue:
<code>legend( ,fill=v)</code>	Colores para rellenar
<code>legend( ,col=v)</code>	Colores de puntos y líneas
<code>legend( ,lty=v)</code>	Tipos de línea
<code>legend( ,lwd=v)</code>	Anchura de línea
<code>legend( ,pch=v)</code>	Caracteres para dibujar (vector de caracteres)
<code>title(main,sub)</code>	Añade un título, <code>main</code> , en la parte superior del gráfico actual, de tamaño grande, y un subtítulo, <code>sub</code> , en la parte inferior, de tamaño menor.

---

<code>axis(side,...)</code>	Añade al gráfico actual un eje en el lado indicado por el primer argumento (de 1 a 4, en el sentido de las agujas del reloj, siendo el 1 la parte inferior). Otros argumentos controlan la posición de los ejes, dentro o fuera del gráfico, las marcas y las etiquetas. Es útil para añadir ejes tras utilizar la función <code>plot()</code> con el argumento <code>axes=F</code> .
-----------------------------	---

---

Las funciones gráficas de nivel bajo necesitan normalmente alguna información de posición, como las coordenadas  $x$  e  $y$ , para determinar dónde colocar los nuevos elementos. Las coordenadas se dan en términos de *coordenadas de usuario* las cuales están definidas por las funciones gráficas de alto nivel previas y se toman en función de los datos suministrados a estas funciones.

Los dos argumentos  $x$  e  $y$ , pueden sustituirse por un solo argumento de clase lista con dos componentes llamados  $x$  e  $y$ , o por una matriz con dos columnas. De este modo, funciones como `locator()`, que vemos a continuación, pueden usarse para especificar interactivamente posiciones en un gráfico.

### 11.2.1 Anotaciones matemáticas

En muchas ocasiones es útil añadir símbolos matemáticos y fórmulas a un gráfico. En R es posible hacerlo especificando una expresión, en vez de una cadena de caracteres en cualquiera de las funciones `text`, `mtext`, `axis` o `title`. Por ejemplo, la orden siguiente dibuja la fórmula de la distribución binomial en la posición  $x, y$ :

```
> text(x, y, expression(paste(bgroup("(", atop(n, x), ")"), p^x, q^n-x)))
```

La información detallada puede obtenerla con las órdenes:

```
> help(plotmath)
> example(plotmath)
```

### 11.2.2 Fuentes vectoriales Hershey

Es posible escribir texto utilizando las fuentes vectoriales Hershey en las funciones `text` y `contour`. Existen tres razones para utilizar estas fuentes:

- Producen mejores resultados, especialmente en pantalla, con textos rotados o de pequeño tamaño.
- Contienen símbolos que pueden no estar disponibles en las fuentes ordinarias, como signos del zodiaco, cartográficos o astronómicos.
- Contienen caracteres cirílicos y japoneses (Kana y Kanji).

La información detallada, incluyendo las tablas de caracteres, puede obtenerla con las órdenes:

```
> help(Hershey)
> example(Hershey)
> help(Japanese)
> example(Japanese)
```

## 11.3 Funciones gráficas interactivas

R dispone de funciones que permiten al usuario extraer o añadir información a un gráfico utilizando el ratón. La más sencilla es la función `locator()`:

---

<code>locator(n,type)</code>	Permite que el usuario seleccione posiciones del gráfico actual, fundamentalmente, utilizando el botón primario del ratón, hasta que haya seleccionado <code>n</code> puntos (el valor predeterminado son 500) o pulse el botón secundario. El argumento <code>type</code> permite dibujar utilizando los puntos seleccionados con el mismo significado que en las funciones de nivel alto. Su valor predeterminado es no dibujar. La función <code>locator()</code> devuelve las coordenadas de los puntos seleccionados en una lista con dos componentes, <code>x</code> e <code>y</code> .
------------------------------	---

---

La función `locator()` suele utilizarse sin argumentos. Es particularmente útil para seleccionar interactivamente posiciones para elementos gráficos tales etiquetas cuando es difícil calcular previamente dónde deben colocarse. Por ejemplo, para colocar un texto informativo junto a un punto anómalo, la orden

```
> text(locator(1), "Anómalo", adj=0)
```

puede ser útil. En el caso de que no se disponga de ratón, la función `locator()` puede ser utilizada; en este caso se preguntarán al usuario las coordenadas de  $x$  e  $y$ .

---

<code>identify(x,y, labels)</code>	Permite identificar cualquiera de los puntos definidos por <code>x</code> e <code>y</code> utilizando el botón primario del ratón, dibujando la correspondiente componente de <code>labels</code> al lado (o el índice del punto si no existe <code>labels</code> ). Al pulsar el botón secundario del ratón, devuelve los índices de los puntos seleccionados.
--	---

---

A veces interesa identificar *puntos* particulares en un gráfico y no sus posiciones. Por ejemplo, si queremos que el usuario seleccione una observación de interés y, posteriormente, manipularla en algún modo. Dado un número de coordenadas,  $(x, y)$ , en dos vectores numéricos, `x` e `y`, podríamos usar la función `identify()` del siguiente modo:

```
> plot(x,y)
```

```
> identify(x,y)
```

La función `identify()` no realiza ningún gráfico por sí misma sino que permite al usuario mover el puntero del ratón y pulsar junto a un punto. El punto más próximo al puntero (si está suficientemente próximo) se identificará dibujando junto a él su número índice, esto es, la posición que ocupa en los vectores `x` e `y`. Alternativamente podría haber utilizado una cadena de caracteres (como por ejemplo el nombre del caso) para la identificación utilizando el argumento `labels`, o inhabilitar la identificación utilizando el argumento `plot=F`. Cuando se pulsa el botón secundario, la función devuelve los índices de los puntos seleccionados, que podrán ser utilizados para obtener los puntos correspondientes de los vectores `x` e `y`.

## 11.4 Uso de parámetros gráficos

Al crear gráficos, especialmente con fines de presentación o publicación, es posible que R no produzca de modo automático exactamente lo que se desea. Ahora bien, es posible personalizar cada aspecto del gráfico utilizando *parámetros gráficos*. R dispone de muchos parámetros gráficos que controlan aspectos tales como estilo de línea, colores, disposición de las figuras y justificación

del texto entre otros muchos. Cada parámetro gráfico tiene un nombre (por ejemplo, 'col', que controla los colores) y un valor (por ejemplo, "blue", para indicar el color azul).

Por cada dispositivo gráfico activo, se mantiene una lista de parámetros gráficos, y cada dispositivo gráfico dispone de un conjunto predeterminado de parámetros cuando se inicializa. Los parámetros gráficos pueden indicarse de dos modos; bien de modo permanente, lo que afectará a todas las funciones gráficas que accedan al dispositivo gráfico, bien temporalmente, lo que sólo afecta a la función gráfica que lo utiliza en ese momento.

#### 11.4.1 Cambios permanentes. La función `par()`

La función `par()` se utiliza para acceder a la lista de parámetros gráficos del dispositivo gráfico actual y para modificarla.

---

<code>par()</code>	Sin argumentos, devuelve una lista de todos los parámetros gráficos y sus valores, para el dispositivo gráfico actual.
<code>par(c("col", "lty"))</code>	Con un vector de caracteres como argumento, devuelve una lista que contiene sólo los valores de los parámetros citados.
<code>par(col=4,lty=2)</code>	Con argumentos con nombre (o una lista como argumento) establece los valores de estos parámetros y devuelve (de modo invisible) una lista con los valores originales de los parámetros.

---

Al modificar algún parámetro con la función `par()`, la modificación es *permanente*, en el sentido de que cualquier llamada a una función gráfica (en el mismo dispositivo gráfico) vendrá afectada por dicho valor. Puede pensarse que este tipo de asignación equivale a modificar los valores predeterminados de los parámetros que utilizarán las funciones gráficas salvo que se les indique un valor alternativo.

Las llamadas a la función `par()` *siempre* afectan a los valores globales de los parámetros gráficos, incluso aunque la llamada se realice desde una función. Esta conducta a menudo no es satisfactoria; habitualmente desearíamos modificar algunos parámetros, realizar algunos gráficos y volver a los valores originales para no afectar a la sesión completa de R. Este trabajo recae en el usuario que debe salvar el resultado de `par()` cuando realice algún cambio y restaurarlos cuando el gráfico esté completo.

```
> par.anterior <- par(col=4,lty=2)
...órdenes gráficas...
> par(par.anterior)
```

#### 11.4.2 Cambios temporales. Argumentos de las funciones gráficas

Los parámetros gráficos también pueden pasarse a prácticamente todas las funciones gráficas como argumentos con nombre, lo que tiene el mismo efecto que utilizarlos en la función `par()`, excepto que los cambios sólo existen durante la llamada a la función. Por ejemplo:

```
> plot(x,y,pch="+")
```

realiza un diagrama de dispersión utilizando el signo de sumar, +, para representar cada punto, sin cambiar el carácter predeterminado para gráficos posteriores.

## 11.5 Parámetros gráficos habituales

A continuación se detallan muchos de los parámetros gráficos habituales. La ayuda de la función `par()` contiene un resumen más conciso; por lo que puede considerar que ésta es una alternativa más detallada.

Los parámetros gráficos se presentarán en la siguiente forma:

---

<code>nombre=valor</code>	Descripción del efecto del parámetro. <b>nombre</b> es el nombre del parámetro, esto es, el nombre del argumento que debe usar en la función <code>par()</code> o cualquier función gráfica. <b>valor</b> es un valor típico del parámetro.
---------------------------	---

---

### 11.5.1 Elementos gráficos

Los gráficos de R están formados por puntos, líneas, texto y polígonos (regiones rellenas). Existen parámetros gráficos que controlan como se dibujan los *elementos gráficos* citados, como los siguientes:

---

<code>pch="+"</code>	Carácter que se utiliza para dibujar un punto. El valor predeterminado varía entre dispositivos gráficos, pero normalmente es 'o'. Los puntos tienden a aparecer en posición levemente distinta de la exacta, salvo que utilice el carácter "." que produce puntos centrados.
<code>pch=4</code>	Si <code>pch</code> toma un valor entre 0 y 18, ambos inclusive, se utiliza un símbolo especial para realizar los gráficos. Para saber cuales son los gráficos, utilice las órdenes  <pre>&gt; plot(1,t="n")</pre> <pre>&gt; legend(locator(1), as.character(0:18), pch=0:18)</pre> y pulse en la parte superior del gráfico.
<code>lty=2</code>	Es el tipo de línea. Algunos tipos no pueden dibujarse en algunos dispositivos gráficos, pero siempre, el tipo 1 corresponde a una línea continua, y los tipos 2 o superior corresponden a líneas con puntos, rayas o combinaciones de ambos.
<code>lwd=2</code>	Es la anchura de línea, medida en múltiplos de la anchura "base". Afecta tanto a los ejes como a las líneas dibujadas con la función <code>lines()</code> , etc.
<code>col=2</code>	Es el color que se utiliza para los puntos, líneas, texto, imágenes y relleno de zonas. Cada uno de estos elementos gráficos admite una lista de colores posibles y el valor de este parámetro es un índice para esta lista. Obviamente este parámetro solo tiene aplicación en algunos dispositivos.
<code>font=2</code>	Es un entero que indica que fuente se utilizará para el texto. Si es posible, 1 corresponde a texto normal, 2 a <b>negrilla</b> , 3 a <i>itálica</i> y 4 a <b><i>itálica negrilla</i></b> .
<code>font.axis,</code> <code>font.lab,</code> <code>font.main,</code> <code>font.sub</code>	Es la fuente que se utilizará, respectivamente, para escribir en los ejes, etiquetado de x e y, título y subtítulo.

---

---

<code>adj=-0.1</code>	Indica cómo debe justificarse el texto respecto de la posición de dibujo. 0 indica justificación a la izquierda, 1 indica justificación a la derecha, 0.5 indica centrado. Puede usar cualquier otro valor que indicará la proporción de texto que aparece a la izquierda de la posición de dibujo, por tanto un valor de -0.1 dejará un 10% de la anchura del texto entre el mismo y la posición de dibujo.
-----------------------	--

---

<code>cex=1.5</code>	Es el carácter de expansión. Su valor indica el tamaño de los caracteres de texto (incluidos los caracteres de dibujo) respecto del tamaño predeterminado.
----------------------	--

---

### 11.5.2 Ejes y marcas de división

Muchos de los gráficos de nivel alto en R tienen ejes, pero, además, siempre es posible construirlos con la función gráfica de nivel bajo `axis()`. Los ejes tienen tres componentes principales: La *línea del eje*, cuyo estilo lo controla el parámetro `lty`, las *marcas de división*, que indican las unidades de división del eje, y las *etiquetas de división*, que indican las unidades de las marcas de división. Estas componentes pueden modificarse con los siguientes parámetros gráficos:

---

<code>lab=c(5,7,12)</code>	Los dos primeros números indican el número de marcas del eje <i>x</i> y del eje <i>y</i> respectivamente. El tercer número es la longitud de las etiquetas de los ejes medida en caracteres (incluido el punto decimal); si elige un número muy pequeño puede conducir a que todas las etiquetas de división se redondeen al mismo número.
----------------------------	--

---

<code>las=1</code>	Corresponde a la orientación de las etiquetas de los ejes. Un 0 indica paralelo al eje, un 1 indica horizontal, y un 2 indica perpendicular al eje.
--------------------	---

---

<code>mgp=c(3,1,0)</code>	Son las posiciones de las componentes de los ejes. La primera es la distancia desde la etiqueta del eje al eje, medida en líneas de texto. La segunda es la distancia hasta las etiquetas de división. La tercera y última es la distancia desde el eje a la línea del eje (normalmente cero). Los valores positivos indican que está fuera de la zona de dibujo, y los negativos indican que está dentro.
---------------------------	--

---

<code>tck=0.01</code>	Es la longitud de las marcas de división, dada como una fracción de la zona de dibujo. Cuando <code>tck</code> es pequeño (menos de 0.5) las marcas de división de ambos ejes serán del mismo tamaño. Un valor de 1 hará que aparezca una rejilla. Si el valor es negativo, las marcas de división se harán por la parte exterior de la zona de dibujo. Utilice <code>tck=0.01</code> y <code>mgp=c(1,-1.5,0)</code> para marcas de división internas.
-----------------------	--

---



`xaxs="s"`  
`yaxs="d"`

Estilo de eje de los ejes  $x$  e  $y$  respectivamente. Con los estilos "**s**" (estándar) y "**e**" (extendido) tanto la mayor marca como la menor caen fuera del intervalo de los datos. Los ejes extendidos pueden ampliarse levemente si hay algún punto muy próximo al borde. Este estilo de ejes puede dejar a veces grandes zonas en blanco cerca de los bordes. Con los estilos "**i**" (interno) y "**r**" (predeterminado) las marcas de división siempre caen dentro del rango de los datos, aunque el estilo "**r**" deja un pequeño espacio en los bordes.

Si selecciona el estilo "**d**" (directo) se procede a *bloquear* los ejes actuales y los utiliza para los siguientes gráficos hasta que el parámetro se cambie a otro de los valores. Este procedimiento es útil para generar series de gráficos con escalas fijas.

### 11.5.3 Márgenes de las figuras

Un gráfico de R se denomina **figura** y comprende una *zona de dibujo* rodeada de márgenes (que posiblemente contendrán etiquetas de ejes, títulos, etc.) y, normalmente, acotada por los propios ejes. Una figura típica aparece en la figura 11. Los parámetros gráficos que controlan la disposición de la figura incluyen:

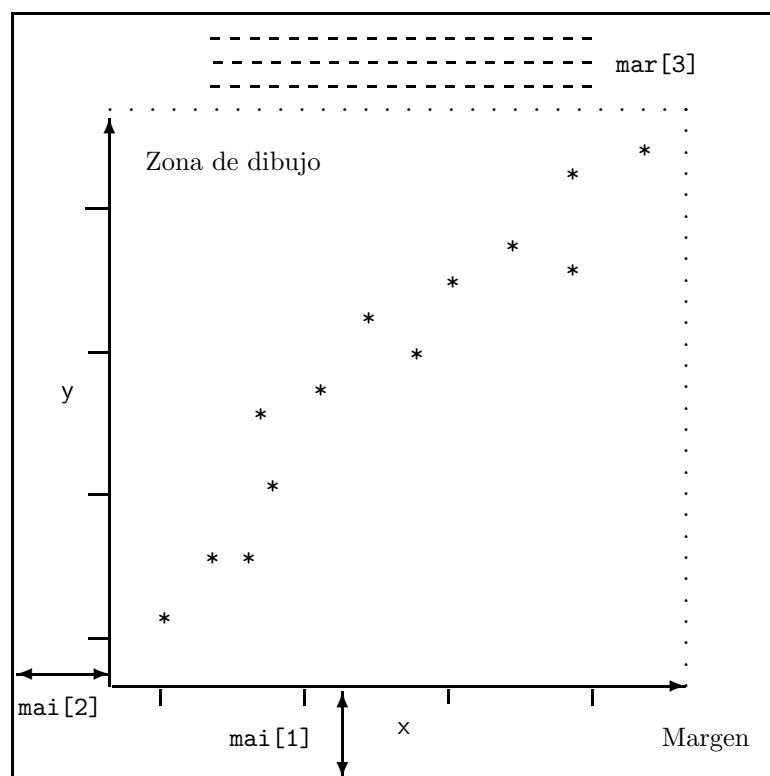


Figura 11: Anatomía de una figura en R

`mai=` Anchuras de los márgenes inferior, izquierdo, superior y derecho respectivamente, medidos en pulgadas.  
`c(1,0.5,0.5,0)`  
`mar=c(4,2,2,1)` Similar a `mai`, pero medido en líneas de texto.

`mar` y `mai` están relacionados en el sentido de que un cambio en uno se refleja en el otro. Los valores predeterminados son a menudo demasiado grandes, el margen derecho se necesita raramente, ni tampoco el superior si no se incluye título. Los márgenes inferior e izquierdo sólo necesitan el tamaño preciso para incluir las etiquetas de ejes y de división. Además el valor predeterminado no tiene en cuenta la superficie del dispositivo gráfico, así, si utiliza el dispositivo `postscript()` con el argumento `height=4` conducirá a un gráfico en que la mitad son márgenes, salvo que explícitamente cambie `mar` o `mai`. Cuando hay figuras múltiples, como veremos después, los márgenes se reducen a la mitad, aunque suele no ser suficiente cuando varias figuras comparten la misma página.

#### 11.5.4 Figuras múltiples

R permite la creación de una matriz de  $n \times m$  figuras en una sola página. Cada figura tiene sus propios márgenes, y la matriz de figuras puede estar opcionalmente rodeada de un *margin exterior*, tal como se muestra en la figura 12.

Los parámetros gráficos relacionados con las figuras múltiples son los siguientes:

<code>mfcol=c(3,2)</code> <code>mfrow=c(2,4)</code>	Definen el tamaño de la matriz de figuras múltiples. En ambos, el primer valor es el número de filas, el segundo el de columnas. La diferencia entre ambas formas, es que con el primero, <code>mfcol</code> , la matriz se rellena por columnas, en tanto que con el segundo, <code>mfrow</code> , lo hace por filas. La distribución en Figure 12 debería crearse con <code>mfrow=c(3,2)</code> y está representada la página cuando se han realizado cuatro gráficos.
<code>mfg=c(2,2,3,2)</code>	Definen la posición de la figura actual dentro de la matriz de figuras múltiples. Los dos primeros valores indican la fila y columna de la figura actual, en tanto que los dos últimos son el número de filas y columnas de la matriz de figuras múltiples. Estos parámetros se utilizan para saltar entre figuras de la matriz. Incluso los dos últimos valores pueden ser distintos de los verdaderos valores, para obtener figuras de tamaños distintos en la misma página.
<code>fig=c(4,9,1,4)/10</code>	Definen la posición de la figura actual en la página. Los valores son las posiciones de los bordes izquierdo, derecho, inferior y superior respectivamente, medidos como la proporción de página desde la esquina inferior izquierda. El ejemplo correspondería a una figura en la parte inferior derecha de la página. Este parámetro permite colocar una figura en cualquier lugar de la página.
<code>oma=c(2,0,3,0)</code> <code>omi=c(0,0,0.8,0)</code>	Definen el tamaño de los márgenes exteriores. De modo similar a <code>mar</code> y <code>mai</code> , el primero está expresado en líneas de texto y el segundo en pulgadas, corresponden a los márgenes inferior, izquierdo, superior y derecho respectivamente.

Los márgenes exteriores son particularmente útiles para ajustar convenientemente los títulos, etc. Puede añadir texto en estos márgenes con la función `mtext()` sin más que utilizar el argumento `outer=T`. R no crea márgenes exteriores predeterminadamente, pero pueden crearse utilizando las funciones `oma` y `omi`.

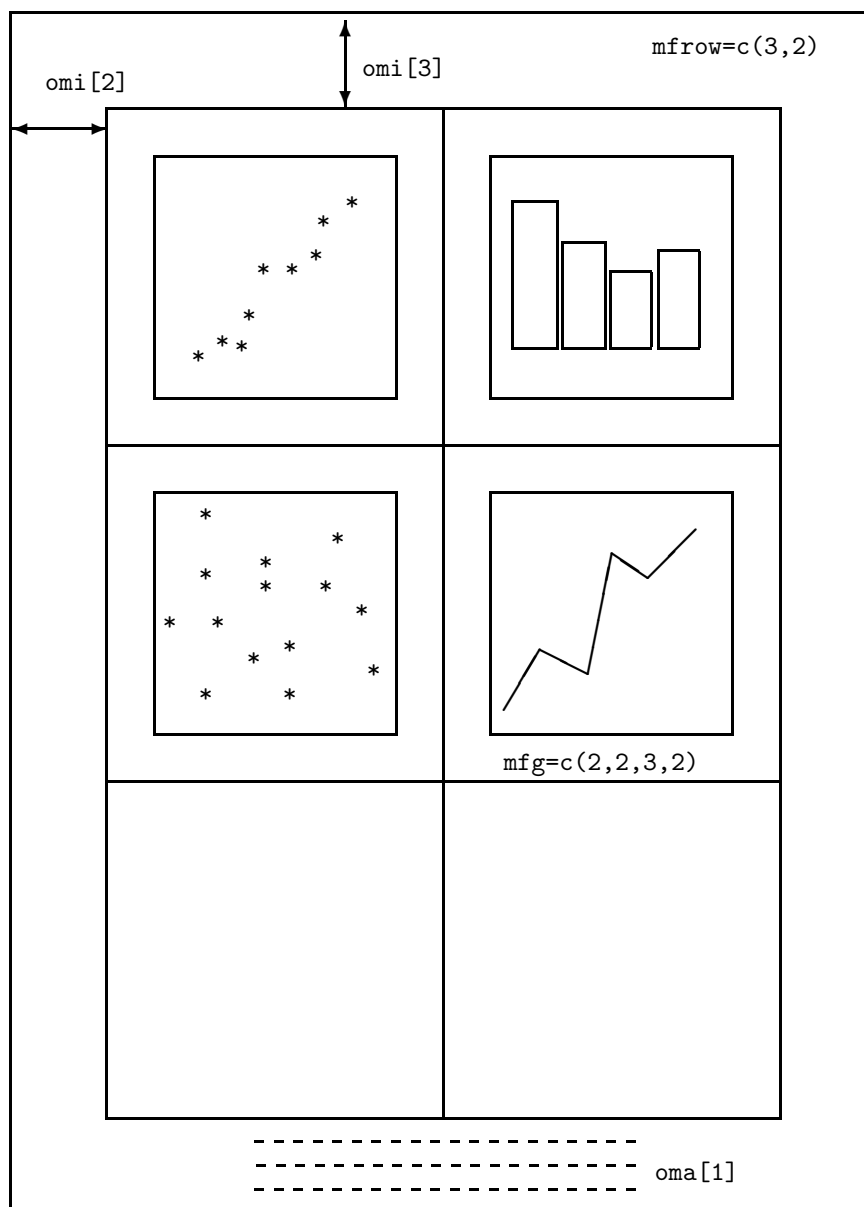


Figura 12: Disposición de página en el modo de figuras múltiples

## 11.6 Dispositivos gráficos

R puede generar gráficos (con niveles de calidad diversos) en casi todos los tipos de pantalla o dispositivos de impresión. Ahora bien, en primer lugar es necesario informarle del tipo de dispositivo de que se trata. Esta acción se realiza iniciando un *controlador de dispositivo*. El propósito de un controlador de dispositivo es convertir las instrucciones gráficas de R (por ejemplo, ‘dibuja una línea’) en una forma que el dispositivo concreto pueda comprender.

Los controladores de dispositivo se inician llamando a una función de controlador de dispositivo. Existe una función para cada controlador de dispositivo y la lista completa puede conseguirse con `help(Devices)`. Por ejemplo, la orden

```
> postscript()
```

dirige cualquier salida gráfica a la impresora en formato PostScript. Algunos controladores de

dispositivo habituales son:

---

<code>X11()</code>	Para uso con los sistemas de ventanas X11 y MICROSOFT WINDOWS.
<code>postscript()</code>	Para imprimir en impresoras PostScript o crear archivos con este formato.
<code>pictex()</code>	Crea un archivo para L <sup>A</sup> T <sub>E</sub> X.

---

Al terminar de utilizar un dispositivo, asegúrese de finalizar el dispositivo utilizando la orden

```
> dev.off()
```

Esta orden asegura que el dispositivo finaliza correctamente; por ejemplo en el caso de una impresora asegura que cada página sea completada y enviada a la impresora.

### 11.6.1 Inclusión de gráficos PostScript en documentos

Si utiliza el argumento `file` en la función `postscript()`, almacenará los gráficos, en formato PostScript, en el archivo que desee. El gráfico tendrá orientación apaisada salvo que especifique el argumento `horizontal=FALSE`. El tamaño del gráfico se controla con los argumentos `width` (anchura) y `height` (altura) que se utilizan como factores de escala para ajustarlo a dichas dimensiones. Por ejemplo, la orden

```
> postscript("archivo.ps", horizontal=FALSE, height=5)
```

producirá un archivo que contiene el código PostScript para una figura de cinco pulgadas de alto, que podrá ser incluido en un documento<sup>26</sup>. Tenga en cuenta que si el archivo ya existe, será previamente borrado. Ello ocurre aunque el archivo se haya creado previamente en la misma sesión de R.

### 11.6.2 Dispositivos gráficos múltiples

La utilización avanzada de R implica a menudo tener varios dispositivos gráficos en uso simultáneamente. Naturalmente solo un dispositivo gráfico acepta las órdenes gráficas en cada momento, y se denomina *dispositivo actual*. Cuando se abren varios dispositivos, forman una sucesión numerada cuyos nombres determinan el tipo de dispositivo en cada posición.

Las principales órdenes utilizadas para trabajar con varios dispositivos y sus significados son las siguientes:

---

<code>X11()</code> <code>postscript()</code> ...	Cada llamada a una función de controlador de dispositivo abre un nuevo dispositivo gráfico y, por tanto, añade un elemento a la lista de dispositivos, al tiempo que este dispositivo pasa a ser el dispositivo actual al que se enviarán los resultados gráficos.
<code>dev.list()</code>	Devuelve el número y nombre de todos los dispositivos activos. El dispositivo de la posición 1 es siempre el <i>dispositivo nulo</i> que no acepta ninguna orden gráfica.
<code>dev.next()</code> <code>dev.prev()</code>	Devuelve el número y nombre del dispositivo gráfico siguiente o anterior, respectivamente, al dispositivo actual.

---

<sup>26</sup> **Advertencia:** El código PostScript producido con esta función *no* es Encapsulated PostScript, y por tanto su inclusión en un documento puede plantear algún problema dependiendo del programa que se utilice. En L<sup>A</sup>T<sub>E</sub>X la inclusión no plantea problemas. Si desea transformar el archivo a formato EPS, puede utilizar GHOSTSCRIPT.

---

<code>dev.set(which=k)</code>	Puede usarse para hacer que el dispositivo actual sea el que ocupa la posición <code>k</code> en la lista de dispositivos. Devuelve el número y nombre del dispositivo.
<hr/>	
<code>dev.off(k)</code>	Cierra el dispositivo gráfico que ocupa la posición <code>k</code> de la lista de dispositivos. Para algunos dispositivos, como los dispositivos <code>postscript</code> , finalizará el gráfico, bien imprimiéndolo inmediatamente, bien completando la escritura en el archivo, dependiendo de cómo se ha iniciado el dispositivo.
<hr/>	
<code>dev.copy(device,</code> <code>..., which=k)</code> <code>dev.print(device,</code> <code>..., which=k)</code>	realiza una copia del dispositivo <code>k</code> . Aquí, <code>device</code> , es una función de dispositivo, como <code>postscript</code> , con argumentos adicionales si es necesario, especificados por <code>...</code> . La función <code>dev.print</code> es similar, pero el dispositivo copiado se cierra inmediatamente, lo que finaliza las acciones pendientes que se realizan inmediatamente. (Vea también <code>printgraph()</code> ).
<hr/>	
<code>graphics.off()</code>	Cierra todos los dispositivos gráficos de la lista, excepto el dispositivo nulo.

---

## A Primera sesión con R

Las siguiente sesión pretende presentar algunos aspectos del entorno R utilizándolos. Muchos de estos aspectos le serán desconocidos e incluso enigmáticos al principio, pero esta sensación desaparecerá rápidamente. La sesión está escrita para un usuario de UNIX. Los usuarios de MACINTOSH y MICROSOFT WINDOWS deben adaptarla apropiadamente.

---

<code>login:</code> <code>...</code>	Conéctese e inicie el sistema de ventanas. Es necesario que copie el archivo <code>morley.data</code> en su directorio de trabajo. Si no sabe hacerlo, consulte con un experto local. Si ya lo ha copiado, continúe.
<code>\$ R</code>	Ejecute R.  Comienza R y aparece el mensaje inicial.  (Dentro de R no mostraremos el símbolo de ‘preparado’ en la parte izquierda para evitar confusiones.)
<code>help.start()</code>	Muestra la ayuda interactiva, que está escrita en formato HTML, utilizando el lector WEB disponible en el ordenador. Si utiliza esta ayuda comprenderá mejor el funcionamiento de R. Minimice la ventana de ayuda y continúe la sesión.
<code>x11()</code>	Abre una ventana gráfica en la pantalla. Es posible que necesita moverla o cambiarla de tamaño para facilitar el trabajo con las dos ventanas.

---

<code>x ← rnorm(50)</code> <code>y ← rnorm(x)</code>	Genera dos vectores que contienen cada uno 50 valores pseudoaleatorios obtenidos de una distribución normal (0,1) y los almacena en <i>x</i> e <i>y</i> .
<code>plot(x, y)</code>	Representa 50 puntos en el plano, tomando cada componente de <i>x</i> y de <i>y</i> como las coordenadas de un punto.
<code>objects()</code>	Presenta los nombres de los objetos existentes en ese momento en la imagen de R.
<code>rm(x,y)</code>	Elimina los objetos <i>x</i> e <i>y</i> .

---

<code>x ← 1:20</code>	Almacena en <i>x</i> el vector (1, 2, ..., 20)
<code>w ← 1 + sqrt(x)/2</code>	A partir del vector <i>x</i> , crea un vector ponderado de desviaciones típicas, y lo almacena en <i>w</i> .
<code>hoja.de.datos ←   data.frame(x=x,           y= x + rnorm(x)*w)</code> <code>hoja.de.datos</code>	Crea una <i>hoja de datos</i> de dos columnas, llamadas <i>x</i> e <i>y</i> , y la almacena en <code>hoja.de.datos</code> . A continuación la presenta en pantalla.
<code>regr ← lm(y~x,           data=hoja.de.datos)</code> <code>summary(regr)</code>	Realiza el ajuste de un modelo de regresión lineal de <i>y</i> sobre <i>x</i> , lo almacena en <code>regr</code> y presenta en pantalla un resumen del análisis.
<code>regr.pon ← lm(y~x,           data=hoja.de.datos,           weight=1/w^2)</code> <code>summary(regr.pon)</code>	Puesto que se conocen las desviaciones típicas, puede realizarse una regresión ponderada.
<code>attach(hoja.de.datos)</code>	Conecta la hoja de datos, de tal modo que sus columnas aparecen como variables.
<code>regr.loc ← lowess(x, y)</code>	Realiza una regresión local no paramétrica.

<code>plot(x, y)</code>	Representa el gráfico de puntos estándar.
<code>lines(x, regr.loc\$y)</code>	Le añade la regresión local.
<code>abline(0, 1, lty=3)</code>	Le añade la verdadera recta de regresión (punto de corte 0, pendiente 1).
<code>abline(coef(regr))</code>	Le añade la recta de regresión no ponderada.
<code>abline(coef(regr.pon), lty=4)</code>	Le añade la recta de regresión ponderada.
<code>detach()</code>	Desconecta la hoja de datos, eliminándola de la lista de búsqueda.
<code>plot(fitted(regr), resid(regr), xlab="Predichos", ylab="Residuos", main= "Residuos / Predichos")</code>	Un gráfico diagnóstico de regresión para investigar la heteroscedasticidad. ¿Puede verla?
<code>qqnorm(resid(regr), main= "Residuos por rangos")</code>	Gráfico en papel probabilístico normal para comprobar asimetría, aplastamiento y datos anómalos. (No es muy útil en este caso)
<code>rm(x,w,hoja.de.datos, regr,regr.pon,regr.loc)</code>	Elimina los objetos creados.

	Ahora analizaremos los datos del experimento clásico de Michaelson y Morley para medir la velocidad de la luz.
<code>system("more morley.data")</code>	Opcional. Interrumpe temporalmente R y presenta el contenido del archivo. La función <code>system</code> permite ejecutar órdenes del sistema operativo.
<code>mm ← read.table( "morley.data") mm</code>	Lee los datos de Michaelson y Morley y los almacena en la hoja de datos <code>mm</code> , y la muestra en pantalla a continuación. Hay cinco experimentos (col. <code>Expt</code> ) y cada uno contiene 20 series (col. <code>Run</code> ) y <code>s1</code> es la velocidad de la luz medida en cada caso, codificada apropiadamente.
<code>mm\$Expt ← factor(mm\$Expt) mm\$Run ← factor(mm\$Run)</code>	Transforma <code>Expt</code> y <code>Run</code> en factores.
<code>attach(mm)</code>	Conecta la hoja de datos en la posición predeterminada: la 2.
<code>plot(Expt,Speed, main= "Velocidad de la luz", xlab="Experimento No.")</code>	Compara los cinco experimentos mediante diagramas de cajas.
<code>fm ← aov(Speed~Run+Expt, data=mm) summary(fm)</code>	Analiza los datos como un diseño en bloques aleatorizados, tomando las 'series' y los 'experimentos' como factores.
<code>fm0 ← update(fm, ~.-Run) anova(fm0,fm)</code>	Ajusta el submodelo eliminando las 'series', y lo compara utilizando un análisis de la varianza.
<code>detach() rm(fm, fm0)</code>	Desconecta la hoja de datos y elimina los objetos creados, antes de seguir adelante.

	Consideraremos ahora nuevas posibilidades gráficas.
<code>x ← seq(-pi, pi, len=50) y ← x</code>	$x$ es un vector con 50 valores equiespaciados en el intervalo $-\pi \leq x \leq \pi$ . $y$ es idéntico.

<code>f ← outer(x, y, function(x,y) cos(y)/(1+x^2))</code>	$f$ es una matriz cuadrada, cuyas filas y columnas están indexadas por $x$ e $y$ respectivamente, formada por los valores de la función $\cos(y)/(1+x^2)$ .
<code>oldpar ← par() par(pty="s")</code>	Almacena los parámetros gráficos y modifica el parámetro <code>pty</code> (zona de dibujo) para que valga "s" (cuadrado).
<code>contour(x, y, f) contour(x, y, f, nlevels=15, add=T)</code>	Dibuja un mapa de curvas de nivel de $f$ ; y le añade más líneas para obtener más detalle.
<code>fa ← (f-t(f))/2</code>	$fa$ es la "parte asimétrica" de $f$ . ( $t(f)$ es la traspuesta de $f$ ).
<code>contour(x, y, fa, nint=15)</code>	Dibuja un mapa de curvas de nivel,...
<code>par(oldpar)</code>	... y recupera los parámetros gráficos originales.
<code>image(x, y, f) image(x, y, fa)</code>	Dibuja dos gráficos de densidad (de los que puede obtener copias impresas si lo desea)
<code>objects(); rm(x,y,f,fa)</code>	Elimina los objetos creados, antes de seguir adelante.
<hr/>	
<code>th ← seq(-pi, pi, len=100) z ← exp(1i*th)</code>	Con R también puede utilizar números complejos. $1i$ se utiliza para representar la unidad imaginaria, $i$
<code>par(pty="s") plot(z, type="l")</code>	La representación gráfica de un argumento complejo consiste en dibujar la parte imaginaria frente a la real. En este caso se obtiene un círculo.
<code>w ← rnorm(100) + rnorm(100)*1i</code>	Suponga que desea generar puntos pseudoaleatorios dentro del círculo unidad. Un primer intento consiste en generar puntos complejos cuyas partes real e imaginaria, respectivamente, procedan de una normal (0,1)...
<code>w ← ifelse(Mod(w) &gt; 1, 1/w, w)</code>	y sustituir los que caen fuera del círculo por sus inversos.
<code>plot(w, xlim=c(-1,1), ylim=c(-1,1), pch="+", xlab="x", ylab="y") lines(z)</code>	Todos los puntos están dentro del círculo unidad, pero la distribución no es uniforme.
<code>w ← sqrt(runif(100))* exp(2*pi*runif(100)*1i) plot(w, xlim=c(-1,1), ylim=c(-1,1), pch="+", xlab="x", ylab="y") lines(z)</code>	Este segundo método utiliza la distribución uniforme. En este caso, los puntos presentan una apariencia más uniformemente espaciada sobre el círculo.
<code>rm(th,w,z)</code>	De nuevo elimina los objetos creados, antes de seguir adelante.
<hr/>	
<code>par(oldpar)</code>	Recupera los parámetros gráficos originales.
<code>butterfly()</code>	El clásico gráfico de mariposa.
<code>rm(oldpar)</code>	Elimina incluso este objeto.
<code>q()</code>	Termina el programa R ...
<code>\$</code>	... y vuelve al sistema operativo.



## B El editor de órdenes de R

### B.1 Preliminares

Si la biblioteca de GNU, **readline**, está disponible cuando se compila R, se puede utilizar un editor interno de líneas de órdenes que permite recuperar, editar y volver a ejecutar las órdenes utilizadas previamente.

El editor puede desactivarse (lo que permite utilizar ESS<sup>27</sup>) dando la opción de inicio ‘--no-readline’.

La versión de MICROSOFT WINDOWS dispone de un editor más sencillo. Vea la opción **Console** en el menú **Help**.

Cuando use R con las posibilidades de **readline**, estarán disponibles las opciones que posteriormente se indican.

Tenga en cuenta las siguientes convenciones tipográficas usuales:

**^M**, significa *Presione la tecla **Control** y, sin soltarla, pulse la tecla **m***.

**Esc m**, significa *Presione la tecla **Esc** y, después de soltarla, pulse la tecla **m***. En este último caso, puede que sea distinto el resultado si se pulsa **m** o **M**.

### B.2 Edición de acciones

R conserva un historial de las órdenes que se teclean, incluyendo las líneas erróneas, lo que permite recuperar las líneas del historial, modificarlas si es necesario, y volver a ejecutarlas como nuevas órdenes. Si el estilo de edición de órdenes es **emacs** cualquier carácter que teclee se inserta en la orden que se está editando, desplazando los caracteres que estén a la derecha del cursor. En el estilo **vi** el modo de inserción de caracteres se inicia con **Esc i** o **Esc a**, y se finaliza con **Esc .**

Cuando pulse la tecla **Return** la orden completa será ejecutada.

Otras acciones posibles de edición se resumen en la tabla siguiente.

### B.3 Resumen del editor de líneas de órdenes

#### 1. Recuperación de órdenes y movimiento vertical

Recupera la orden anterior (retrocede en el historial)    **^P**

Recupera la orden posterior (avanza en el historial)    **^N**

Recupera la última orden que contenga la cadena **texto** **^R texto**

#### 2. Movimiento horizontal del cursor

Va al principio de la línea    **^A**

Va al final de la línea    **^E**

Retrocede una palabra    **Esc b**

Avanza una palabra    **Esc f**

Retrocede un carácter    **^B**

---

<sup>27</sup>Corresponde al acrónimo del editor de textos, ‘Emacs Speaks Statistics’; vea la dirección <http://ess.stat.wisc.edu/>

Avanza un carácter	<code>^F</code>
--------------------	-----------------

### 3. Edición

Inserta <b>texto</b> en el cursor	<code>texto</code>
Añade <b>texto</b> tras el cursor	<code>^Ftexto</code>
Borra el carácter a la izquierda del cursor	<code>Delete</code>
Borra el carácter bajo el cursor	<code>^D</code>
Borra el resto de la palabra bajo el cursor, y la guarda	<code>Esc d</code>
Borra el resto de la línea desde el cursor, y lo guarda	<code>^K</code>
Inserta el último texto guardado	<code>^Y</code>
Intercambia el carácter bajo el cursor con el siguiente	<code>^T</code>
Cambia el resto de la palabra a minúsculas	<code>Esc l</code>
Cambia el resto de la palabra a mayúsculas	<code>Esc c</code>
Vuelve a ejecutar la línea	<code>Return</code>

Al pulsar **Return** termina la edición de la línea.

## Referencias

- D. M. Bates and D. G. Watts (1988), *Nonlinear Regression Analysis and Its Applications*. John Wiley & Sons, New York.
- Richard A. Becker, John M. Chambers and Allan R. Wilks (1988), *The New S Language*. Chapman & Hall, New York. Este libro se conoce como el *Libro azul*.
- John M. Chambers and Trevor J. Hastie eds. (1992), *Statistical Models in S*. Chapman & Hall, New York. Este libro se conoce como el “*Libro blanco*”.
- Annette J. Dobson (1990), *An Introduction to Generalized Linear Models*, Chapman and Hall, London.
- Peter McCullagh and John A. Nelder (1989), *Generalized Linear Models*. Second edition, Chapman and Hall, London.
- John A. Rice (1995), *Mathematical Statistics and Data Analysis*. Second edition. Duxbury Press, Belmont, CA.
- S. D. Silvey (1970), *Statistical Inference*. Penguin, London.