

# Importar Datos

with readr, tibble, and tidyr

Guía Rápida



El **tidyverse** de R está construido alrededor de los **datos ordenados** almacenados en **tibbles**, una versión mejorada del data frame.



El anverso de esta guía muestra como leer ficheros de texto en R con **readr**.



El reverso muestra como crear tibbles con **tibble** y como disponer datos ordenados con **tidyr**.

## Otros tipos de datos

Prueba uno de los siguientes paquetes para importar otros tipos de ficheros

- **haven** - ficheros SPSS, Stata y SAS
- **readxl** - ficheros Excel (.xls y .xlsx)
- **DBI** - bases de datos
- **jsonlite** - json
- **xml2** - XML
- **httr** - Web APIs
- **rvest** - HTML (Web Scraping)

## Funciones Escritura

Salva **x**, un objeto R, en **path**, una ruta del fichero, con:

**write\_csv**(x, path, na = "NA", append = FALSE, col\_names = !append)

**Tibble/df a fichero delimitado con coma.**

**write\_delim**(x, ruta, delim = " ", na = "NA", append = FALSE, col\_names = !append)

**Tibble/df to file with any delimiter.**

**write\_excel\_csv**(x, path, na = "NA", append = FALSE, col\_names = !append)

**Tibble/df a CSV para Excel**

**write\_file**(x, path, append = FALSE)

**Cadena a fichero.**

**write\_lines**(x, path, na = "NA", append = FALSE)

**Vector cadena a fichero, un elemento por línea.**

**write\_rds**(x, path, compress = c("none", "gz", "bz2", "xz"), ...)

**Objeto a fichero RDS.**

**write\_tsv**(x, path, na = "NA", append = FALSE, col\_names = !append)

**Tibble/df a ficheros delimitados por tab.**

## Funciones Lectura

### Lectura de datos tabulares a tibbles

Estas funciones comparten los siguientes argumentos comunes:

```
read_* (file, col_names = TRUE, col_types = NULL, locale = default_locale(), na = c("", "NA"),
        quoted_na = TRUE, comment = "", trim_ws = TRUE, skip = 0, n_max = Inf, guess_max =
        min(1000, n_max), progress = interactive())
```

a,b,c  
1,2,3  
4,5,NA

A	B	C
1	2	3
4	5	NA

**read\_csv()**

Lee ficheros delimitados por comas.

**read\_csv("file.csv")**

a;b;c  
1;2;3  
4;5;NA

A	B	C
1	2	3
4	5	NA

**read\_csv2()**

Lee ficheros delimitados por punto y coma.

**read\_csv2("file2.csv")**

a|b|c  
1|2|3  
4|5|NA

A	B	C
1	2	3
4	5	NA

**read\_delim**(delim, quote = "\"", escape\_backslash = FALSE, escape\_double = TRUE)

Lee ficheros con cualquier delimitador.

**read\_delim("file.txt", delim = "|")**

a b c  
1 2 3  
4 5 NA

A	B	C
1	2	3
4	5	NA

**read\_fwf**(col\_positions)

Lee ficheros con ancho fijo.

**read\_fwf("file.fwf", col\_positions = c(1, 3, 5))**

**read\_tsv()**

Lee ficheros delimitados por tab. También **read\_table()**.

**read\_tsv("file.tsv")**

### Argumentos útiles

a,b,c  
1,2,3  
4,5,NA

**Fichero ejemplo**

**write\_csv**(path = "file.csv",  
x = read\_csv("a,b,c\n1,2,3\n4,5,NA"))

1	2	3
4	5	NA

**Salta líneas**

**read\_csv**("file.csv",  
skip = 1)

A	B	C
1	2	3
4	5	NA

**Sin cabecera**

**read\_csv**("file.csv",  
col\_names = FALSE)

A	B	C
1	2	3

**Lee un subconjunto**

**read\_csv**("file.csv",  
n\_max = 1)

x	y	z
A	B	C
1	2	3
4	5	NA

**Proporciona cabecera**

**read\_csv**("file.csv",  
col\_names = c("x", "y", "z"))

A	B	C
1	2	3
NA	NA	NA

**Valores Faltantes**

**read\_csv**("file.csv",  
na = c("4", "5", ""))

### Lectura de datos no tabulares

**read\_file**(file, locale = default\_locale())

Lee un fichero en una sola cadena.

**read\_file\_raw**(file)

Lee un fichero en un vector raw.

**read\_lines**(file, skip = 0, n\_max = -1L, locale = default\_locale(), na = character(), progress = interactive())

Lee cada línea en una cadena.

**read\_lines\_raw**(file, skip = 0, n\_max = -1L, progress = interactive())

Lee cada línea en un vector raw.

**read\_log**(file, col\_names = FALSE, col\_types = NULL, skip = 0, n\_max = -1, progress = interactive())

Ficheros de log estilo Apache.

## Filtrando tipos de datos

Las funciones de readr interpretan los tipos de cada columna y convierten los tipos de forma apropiada (pero NO convertirán cadenas a factores automáticamente).

Un mensaje muestra el tipo de columna en el resultado.

```
## Parsed with column specification:
## cols(
##   age = col_integer(),
##   sex = col_character(),
##   earn = col_double()
## )
#> # A tibble: 1 x 3
#>   age sex   earn
#>   <int> <chr> <dbl>
#> 1   41 M     45420
```

age es un entero

sex es un caracter

1. Usa **problems()** para diagnosticar problemas  
**x <- read\_csv("file.csv"); problems(x)**

2. Usa **col\_** function para guiar el filtrado

- **col\_guess()** - por defecto
- **col\_character()**
- **col\_double()**
- **col\_euro\_double()**
- **col\_datetime**(format = "") También **col\_date**(format = "") and **col\_time**(format = "")
- **col\_factor**(levels, ordered = FALSE)
- **col\_integer()**
- **col\_logical()**
- **col\_number()**
- **col\_numeric()**
- **col\_skip()**

**x <- read\_csv("file.csv", col\_types = cols(
 A = col\_double(),
 B = col\_logical(),
 C = col\_factor()
))**

3. Para el resto, los lee como vectores carácter y luego los filtra con la función **parse\_**.

- **parse\_guess**(x, na = c("", "NA"), locale = default\_locale())
- **parse\_character**(x, na = c("", "NA"), locale = default\_locale())
- **parse\_datetime**(x, format = "", na = c("", "NA"), locale = default\_locale()) Also **parse\_date()** and **parse\_time()**
- **parse\_double**(x, na = c("", "NA"), locale = default\_locale())
- **parse\_factor**(x, levels, ordered = FALSE, na = c("", "NA"), locale = default\_locale())
- **parse\_integer**(x, na = c("", "NA"), locale = default\_locale())
- **parse\_logical**(x, na = c("", "NA"), locale = default\_locale())
- **parse\_number**(x, na = c("", "NA"), locale = default\_locale())

**x\$A <- parse\_number(x\$A)**

## Tibbles - un data frame mejorado

El paquete **tibble** proporciona una nueva clase S3 para almacenar datos tabulares, el tibble. Tibbles heredan la clase del data frame, pero mejora dos comportamientos:

- **Visualiza** - Cuando se imprime un tibble, R proporciona una vista concisa de los datos que se ajustan en una pantalla.
- **Subconjunto** - [ siempre devuelve un nuevo tibble, [[ y \$ siempre devuelve un vector.
- **No emparejado parcial** - Se debe usar el nombre completo de las columnas cuando se selecciona un subconjunto.

```
# A tibble: 234 x 6
  manufacturer <chr> model <chr> displ <dbl>
1 audi a4 1.8
2 audi a4 1.8
3 audi a4 2.0
4 audi a4 2.0
5 audi a4 2.8
6 audi a4 2.8
7 audi a4 3.1
8 audi a4 quattro 1.8
9 audi a4 quattro 1.8
10 audi a4 quattro 2.0
# ... with 224 more rows, and 3
# more variables: year <int>,
# cyl <int>, trans <chr>
```

### visualización tibble

```
156 1999 6 auto(l4)
157 1999 6 auto(l4)
158 2008 6 auto(l4)
159 2008 8 auto(s4)
160 1999 4 manual(m5)
161 1999 4 auto(l4)
162 2008 4 manual(m5)
163 2008 4 manual(m5)
164 2008 4 auto(l4)
165 2008 4 auto(l4)
166 1999 4 auto(l4)
# reached getOption("max.print")
# -- omitted 68 rows --
```

### Vista de una tabla larga

### visualización data frame

- La apariencia por defecto se controla con las opciones:  
`options(tibble.print_max = n, tibble.print_min = m, tibble.width = Inf)`
- La vista del conjunto completo de datos con **View(x, title)** o **glimpse(x, width = NULL, ...)**
- Revertir a data frame con **as.data.frame()** (requerido por algunos paquetes antiguos)

## Construir un tibble en dos formas

**tibble(...)**  
Construye por columnas.

**tibble(x = 1:3,  
y = c("a", "b", "c"))**

Ambos crean este tibble

**tribble(...)**  
Construye por filas.

**tribble(  
~x, ~y,  
1, "a",  
2, "b",  
3, "c")**

```
A tibble: 3 x 2
  x     y
<int> <dbl>
1     1  a
2     2  b
3     3  c
```

**as\_tibble(x, ...)** Convierte data frame a tibble.

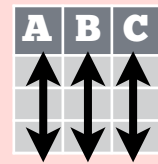
**enframe(x, name = "name", value = "value")**  
Convierte un vector con nombre a un tibble con una columna con nombre y una columna valor.

**is\_tibble(x)** Comprueba si x es un tibble.

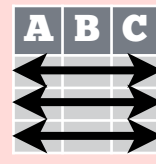
## Datos Ordenados con tidyr

**Datos Tidy** es una forma de organizar datos tabulares. Proporciona una estructura consistente de datos entre paquetes.

Una tabla es tidy si:

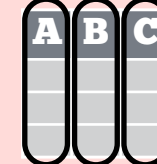


Cada **variable** está en su propia **columna**

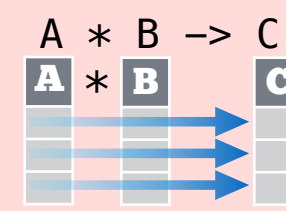


Cada **observación** o **caso**, está en su propia **fila**

Datos Tidy:



Permite el acceso a las variables como vectores



Preserva los casos en las operaciones vectorizadas

## Remodelado de Datos - cambia la disposición de los valores en una tabla

Usa **gather()** y **spread()** para reorganizar los valores de una tabla en una nueva disposición. Usan la idea de una columna clave: par valor columna.

**gather(data, key, value, ..., na.rm = FALSE, convert = FALSE, factor\_key = FALSE)**

Gather mueve los nombres de las columnas a una columna llave, reuniendo los valores de las columnas en una única columna.

table4a

country	1999	2000
A	0.7K	2K
B	37K	80K
C	212K	213K



country	year	cases
A	1999	0.7K
B	1999	37K
C	1999	212K
A	2000	2K
B	2000	80K
C	2000	213K

llave valor

**gather(table4a, `1999`, `2000`,  
key = "year", value = "cases")**

**spread(data, key, value, fill = NA, convert = FALSE, drop = TRUE, sep = NULL)**

Spread mueve el valor único de una columna llave a nombres de columna, repartiendo los valores de una columna entre las nuevas columnas que resultan.

table2

country	year	type	count
A	1999	cases	0.7K
A	1999	pop	19M
A	2000	cases	2K
A	2000	pop	20M
B	1999	cases	37K
B	1999	pop	172M
B	2000	cases	80K
B	2000	pop	174M
C	1999	cases	212K
C	1999	pop	1T
C	2000	cases	213K
C	2000	pop	1T

llave valor

**spread(table2, type, count)**

country	year	cases	pop
A	1999	0.7K	19M
A	2000	2K	20M
B	1999	37K	172M
B	2000	80K	174M
C	1999	212K	1T
C	2000	213K	1T

## Gestión de Datos Faltantes

**drop\_na(data, ...)**

Elimina columnas con NAs.

x	x1	x2
A	1	1
B	NA	NA
C	NA	NA
D	3	3
E	NA	NA

**drop\_na(x, x2)**

**fill(data, ..., .direction = c("down", "up"))**

Completa los NA's en ... columnas con los valores no-NA más cercanos.

x	x1	x2
A	1	1
B	NA	1
C	NA	1
D	3	3
E	NA	3

**fill(x, x2)**

**replace\_na(data, replace = list(), ...)**

Reemplaza NA's por columna.

x	x1	x2
A	1	1
B	NA	2
C	NA	2
D	3	3
E	NA	2

**replace\_na(x, list(x2 = 2), x2)**

## Expansión de Tablas - crea tablas rápidamente con combinaciones de valores

**complete(data, ..., fill = list())**

Añade a los datos combinaciones faltantes de los valores listados en ...

**complete(mtcars, cyl, gear, carb)**

**expand(data, ...)**

Crea un nuevo tibble con todas las posibles combinaciones de valores de las variables listadas en ...

**expand(mtcars, cyl, gear, carb)**

## Separar y Combinar Celdas

Usa estas funciones para separar o combinar celdas en valores individuales, aislados.

**separate(data, col, into, sep = "[^:alnum:]]+", remove = TRUE, convert = FALSE, extra = "warn", fill = "warn", ...)**

Separa cada celda de una columna para crear varias columnas

table3

country	year	rate
A	1999	0.7K/19M
A	2000	2K/20M
B	1999	37K/172M
B	2000	80K/174M
C	1999	212K/1T
C	2000	213K/1T



country	year	cases	pop
A	1999	0.7K	19M
A	2000	2K	20M
B	1999	37K	172
B	2000	80K	174
C	1999	212K	1T
C	2000	213K	1T

**separate\_rows(table3, rate,  
into = c("cases", "pop"))**

**separate\_rows(data, ..., sep = "[^:alnum:]]+", convert = FALSE)**

Separa cada celda en una columna para crear varias filas. También **separate\_rows\_()**.

table3

country	year	rate
A	1999	0.7K/19M
A	2000	2K/20M
B	1999	37K/172M
B	2000	80K/174M
C	1999	212K/1T
C	2000	213K/1T



country	year	rate
A	1999	0.7K
A	1999	19M
A	2000	2K
A	2000	20M
B	1999	37K
B	1999	172M
B	2000	80K
B	2000	174M
C	1999	212K
C	1999	1T
C	2000	213K
C	2000	1T

**separate\_rows(table3, rate)**

**unite(data, col, ..., sep = "\_", remove = TRUE)**

Collapsa celdas de varias columnas para crear una única columna.

table5

country	century	year
Afghan	19	99
Afghan	20	0
Brazil	19	99
Brazil	20	0
China	19	99
China	20	0



country	year
Afghan	1999
Afghan	2000
Brazil	1999
Brazil	2000
China	1999
China	2000

**unite(table5, century, year,  
col = "year", sep = "")**