

# Importar Datos: : GUÍA RÁPIDA



El **tidyverse** de R se basa en **datos ordenados** (tidy data) almacenados en **tibbles**, que son data frames mejorados.



El frente de esta hoja muestra cómo leer archivos de texto en R con **readr**.



El reverso muestra cómo crear tibbles con **tibble** y diseñar datos ordenados con **tidyr**.

## OTROS TIPOS DE DATOS

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

- **haven** – archivos SPSS, Stata y SAS
- **readxl** – archivos excel (.xls y .xlsx)
- **DBI** – base de datos
- **jsonlite** – json
- **xml2** – XML
- **httr** – Web APIs
- **rvest** – HTML (Web Scraping)

## Guardar datos

Guardar **x**, un objeto de R, a **path**, una ruta de acceso a un archivo, como:

### Archivo separado por comas

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

### Archivo con separador arbitrario

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

### CSV para excel

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

### Cadena (string) a archivo

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

### Vector de cadena a archivo, un elemento por línea

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

### Objeto a archivo RDS

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

### Archivo separado por tabulaciones

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

## Leer datos tabulares - Estas funciones comparten estos argumentos:

**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

### Archivo separado por comas

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

Para generar archivo .csv ejecuta:

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

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

A	B	C
1	2	3
4	5	NA

### Archivo separado por punto y coma

**read\_csv2**("archivo2.csv")

**write\_file**(x = "a;b;c\n1;2;3\n4;5;NA", path = "archivo 2.csv")

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

A	B	C
1	2	3
4	5	NA

### Archivo con cualquier separador

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

**write\_file**(x = "a|b|c\n1|2|3\n4|5|NA", path = "archivo.txt")

a b c  
1 2 3  
4 5 NA

A	B	C
1	2	3
4	5	NA

### Archivos de ancho fijo

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

**write\_file**(x = "a b c\n1 2 3\n4 5 NA", path = "archivo.fwf")

### Archivo separado por tabulaciones

**read\_tsv**("archivo.tsv") también con **read\_table()**.

**write\_file**(x = "a\tb\tc\n1\t2\t3\n4\t5\tNA", path = "archivo.tsv")

## ARGUMENTOS ÚTILES

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

### Archivo de ejemplo

**write\_file**("a,b,c\n1,2,3\n4,5,NA", "archivo.csv")

f <- "archivo.csv"

A	B	C
1	2	3
4	5	NA

### Sin encabezado

**read\_csv**(f, col\_names = FALSE)

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

### Proporcionar el encabezado

**read\_csv**(f, col\_names = c("x", "y", "z"))

1	2	3
4	5	NA

### Saltar líneas

**read\_csv**(f, skip = 1)

A	B	C
1	2	3

### Leer un subconjunto

**read\_csv**(f, n\_max = 1)

A	B	C
NA	2	3
4	5	NA

### Valores faltantes

**read\_csv**(f, na = c("1", "."))

## Leer datos no tabulares

### Leer un archivo en una sola cadena

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

### Lee cada línea en una cadena

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

### Lee ficheros de log estilo Apache

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

### Leer un archivo en un vector

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

### Lee cada línea en un vector

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

## Tipos de datos

Las funciones de readr interpretan los tipos de cada columna y convierten los tipos cuando corresponde (pero NO convertirá cadenas en factores automáticamente).

Un mensaje muestra el tipo de cada columna en el resultado.

## Con especificación de columnas:

```
## cols(  
##   age = col_integer(),  
##   sex = col_character(),  
##   earn = col_double()  
## )
```

age es  
un entero

sex es  
un caracter

earn es un decimal  
(numérico)

1. Usa **problems()** para diagnosticar problemas.

**x <- read\_csv("archivo.csv"); problems(x)**

2. Usa una **col\_** function para guiar el parseado.

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

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

3. Sino, leé como vectores de caracteres y luego parsea con una **parse\_** function.

- **parse\_guess()**
- **parse\_character()**
- **parse\_datetime()** También **parse\_date()** y **parse\_time()**
- **parse\_double()**
- **parse\_factor()**
- **parse\_integer()**
- **parse\_logical()**
- **parse\_number()**

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

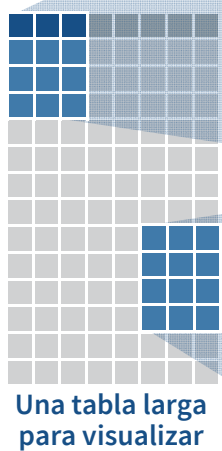


## Tibbles - un data frame mejorado

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



- **Creación de subconjuntos** - [ siempre retorna un nuevo tibble, [[ y \$ siempre retornan un vector.
- **Sin coincidencias parciales**- se deben utilizar los nombres completos de las columnas al crear subconjuntos
- **Impresión en la consola** - Cuando imprimes un tibble, R proporciona una vista concisa de los datos que se ajustan a una pantalla



```
# A tibble: 234 x 6
#   manufacturer model      displ
#   <chr>         <chr>     <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>
```

Vista 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 data frame

- Controla la apariencia por defecto con las opciones:  
`options(tibble.print_max = n,  
tibble.print_min = m, tibble.width = Inf)`
- Ver el conjunto de datos completo con **View()** o **glimpse()**
- Convierte a data frame con **as.data.frame()**

### CONSTRUYE UN TIBBLE DE DOS MANERAS

**tibble(...)**  
Construir por columnas.  
**tibble(x = 1:3, y = c("a", "b", "c"))**

Ambos crean este tibble

**tribble(...)**  
Construir por filas.  
**tribble(~x, ~y,  
1, "a",  
2, "b",  
3, "c")**

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

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

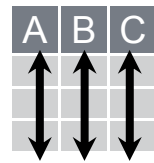
**enframe(x, name = "name", value = "value")**  
Convierte un vector con nombre en un tibble

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

## Tidy Data con tidyr

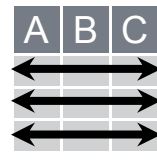
**Tidy data** es una forma de organizar datos tabulares. Proporciona una estructura de datos consistente entre paquetes

Una tabla es tidy si:



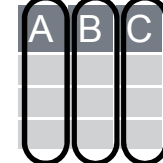
Cada **variable** es una **columna**

&

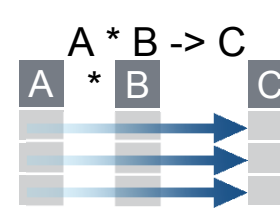


Cada **observación** o **caso** es una **fila**

Tidy data:



Facilita el acceso a las variables como vectores.



Preserva los casos durante las operaciones vectorizadas.

## Remoldear Datos - cambiar la forma de los valores en una tabla

Usa **gather()** y **spread()** para reorganizar los valores de una tabla en una nueva forma.

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

**gather()** mueve los nombres de las columnas a una columna **key**, reuniendo los valores de la columna en una sola columna **value**.

table4a

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

key value

`gather(table4a, `1999`, `2000`,  
key = "anio", value = "casos")`

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

**spread()** mueve los valores únicos de una columna **key** como nombres de columnas, esparciendo los valores de una columna **value** a través de las nuevas columnas.

table2

pais	anio	tipo	cuenta
A	1999	casos	0.7K
A	1999	pop	19M
A	2000	casos	2K
A	2000	pop	20M
B	1999	casos	37K
B	1999	pop	172M
B	2000	casos	80K
B	2000	pop	174M
C	1999	casos	212K
C	1999	pop	1T
C	2000	casos	213K
C	2000	pop	1T

key value

`spread(table2, tipo, cuenta)`

## Manejar Valores Perdidos

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

Elimina las filas que contienen NA's en las columnas especificadas (...)

x

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

drop\_na(x, x2)

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

Completa los NA's en las columnas especificadas (...) con el valor no-NA más cercano.

x

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

fill(x, x2)

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

Reemplaza NA's por columna.

x

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

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

## Expandir tablas - crea rápidamente tablas con combinaciones de valores

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

Completa los datos con combinaciones faltantes de las variables listadas en ... usando los valores de fill o NA's.

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

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

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

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

## Separar Celdas



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

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

Separa cada celda en una columna para hacer varias columnas.

table3

pais	anio	tasa
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

separate(table3, tasa, sep = "/", into = c("casos", "pob"))

pais	anio	casos	pob
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

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

Separa cada celda en una columna para formar varias filas.

table3

pais	anio	tasa
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

separate\_rows(table3, tasa, sep = "/")

pais	anio	tasa
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

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

Une múltiples columnas en una única columna

table5

pais	siglo	anio
Afghan	19	99
Afghan	20	00
Brazil	19	99
Brazil	20	00
China	19	99
China	20	00

unite(table5, siglo, anio, col = "anio", sep = "")

pais	anio
Afghan	1999
Afghan	2000
Brazil	1999
Brazil	2000
China	1999
China	2000

