

readr

Entornos de Análisis de Datos: R

Alberto Torres Barrán

2019-12-10

readr

Introducción

- Paquete para importar y exportar ficheros de texto
- Importar datos:
 - `read_csv()` , para ficheros CSV
 - `read_csv2()` , para ficheros CSV separados por ";"
 - `read_delim()` , para ficheros ASCII delimitados por otros caracteres distintos de "," y ";"
 - `read_tsv()` , para ficheros ASCII delimitados por tabuladores
 - `read_table()` , para ficheros ASCII delimitados por espacios
 - `read_rds()` , formato binario específico de R
- Exportar datos: `write_csv()` . `write_csv2()` , etc.

Ejemplo

```
write_csv(mpg, "mpg.csv")
mpg1 <- read_csv("mpg.csv")
head(mpg1)
```

```
## # A tibble: 6 x 11
##   manufacturer model displ  year   cyl trans      drv   cty   hwy fl   class
##   <chr>         <chr> <dbl> <dbl> <dbl> <chr>   <chr> <dbl> <dbl> <chr> <chr>
## 1 audi         a4      1.8  1999     4 auto(l5) f      18    29 p    compact
## 2 audi         a4      1.8  1999     4 manual(m5) f      21    29 p    compact
## 3 audi         a4      2    2008     4 manual(m6) f      20    31 p    compact
## 4 audi         a4      2    2008     4 auto(av) f      21    30 p    compact
## 5 audi         a4      2.8  1999     6 auto(l5) f      16    26 p    compact
## 6 audi         a4      2.8  1999     6 manual(m5) f      18    26 p    compact
```

Directorio de trabajo

- Directorio donde apunta RStudio

```
getwd()
```

```
## [1] "C:/Users/alberto/Desktop/curso-uah-eadr/src"
```

- Se puede cambiar con `setwd()` o en la pestaña `Files` de RStudio
- Directorio por defecto donde se buscan los ficheros a importar
- Alternativamente, podemos especificar el path completo o usar la herramienta gráfica de RStudio

Missing values en R

- `NA` es una constante que representa valores que faltan (*missing values*)
- Puede estar contenida dentro de vectores (columnas) de cualquier tipo
- `is.na()` devuelve `TRUE` si el valor es `NA` y `FALSE` en caso contrario
- Muchas funciones de R tienen un parámetro opcional `na.rm` que ignora `NA`s

```
dia <-  
  diamonds %>%  
    mutate(y = ifelse(!between(y, 3, 20), NA, y))  
  
dia %>%  
  summarize(y_na = sum(is.na(y)))
```

```
## # A tibble: 1 x 1  
##   y_na  
##   <int>  
## 1     9
```

```
dia %>%  
  summarize(avg_y = mean(y))
```

```
## # A tibble: 1 x 1  
##   avg_y  
##   <dbl>  
## 1    NA
```

```
dia %>%  
  summarize(avg_y = mean(y, na.rm = TRUE))
```

```
## # A tibble: 1 x 1  
##   avg_y  
##   <dbl>  
## 1  5.73
```

Parámetros opcionales

- `col_names` , si TRUE, la primera fila es el nombre de las variables. También se le puede pasar un vector de cadenas de caracteres con los nombres.
- `delim` , carácter que separa las columnas (solo en `read_delim()`)
- `na` , vector con cadenas que se interpretan como missing values. Por defecto `NA` y la cadena vacía.
- `col_types` , vector de clases para las columnas (ver documentación de `col()`). Por defecto se intenta adivinar el tipo de cada columna a partir de las 1000 primeras líneas.
- `n_max` , número máximo de líneas a leer del fichero
- `skip` , número de líneas a ignorar al principio del fichero.
- `locale` , parámetro que nos permite cambiar el encoding, separador decimal y formato de fechas (ver documentación de `locale()`)
- `comment` , una cadena de caracteres que identifica líneas de texto a ignorar (comentarios)
- `trim_ws` , si vale TRUE, se eliminan los espacios en blanco al principio y al final de cada campo

Otros formatos

- `readr` solo tiene funciones para importar ficheros de texto
- Para otros formatos, existen librerías específicas:
 - `haven`, para ficheros de SPSS, Stata y SaS
 - `readxl`, para ficheros de Excel
 - `DBI` junto con otro paquete específico dependiendo de la BD (`RMySQL`, `RSQLite`, etc.) nos permite hacer *queries* contra una BD
 - `jsonlite`, para ficheros JSON
 - `xml2`, para ficheros XML