

Taller 2: Exploración y validación de datos

José Daniel Conejeros

Magíster(c) Sociología UC - Licenciado en Ciencias Sociales

Junio 2020

 Código del taller  jdconejeros@uc.cl  JDConejeros

 [joseconejerosp](#)  [Jose_Conejeros](#)

¿Qué hemos visto hasta el momento?



- Instalación de R y RStudio
- Explorar espacios de trabajo
- Comprender la lógica de objetos
- Explorar vectores, listas y matrices
- Estructura de funciones y argumentos
- Instalar y cargar librerías/paquetes

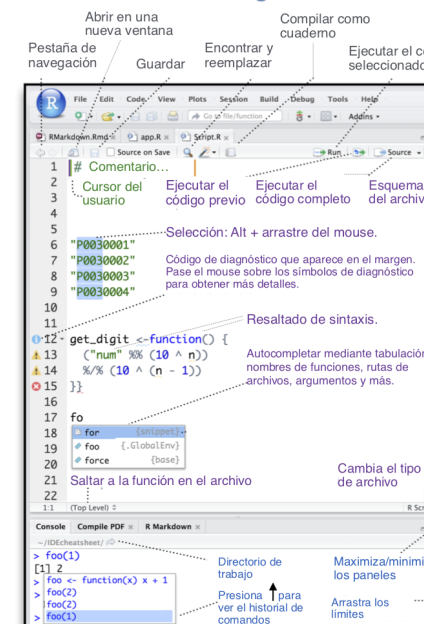
Documentos y apps

RStudio desarrolla herramientas gratuitas y abiertas para R. Su entorno de desarrollo integrado (IDE) facilita el análisis de datos con R.

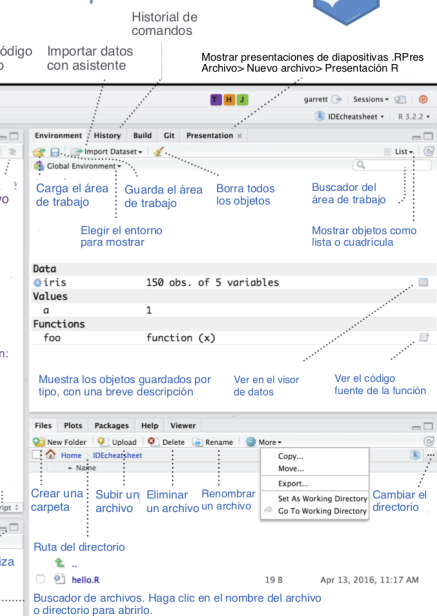
También ofrece: muchos paquetes R (e.g. tidyverse, sparklyr, ggplot2, dplyr), incluidos **Shiny** (crea aplicaciones web sencillas en R) y **R Markdown** (te permite convertir tus análisis en documentos, informes, presentaciones y paneles de alta calidad; compartílos y reproducílos).



Escribe tu código



Soporte R

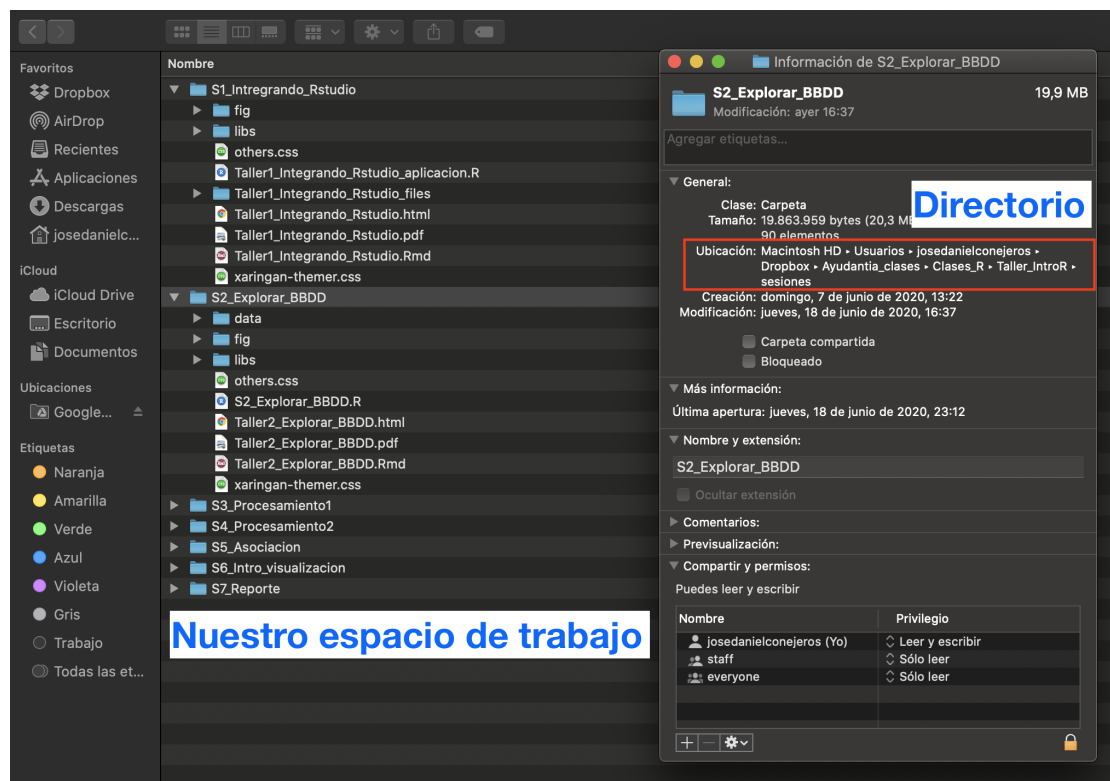


¿Alguna pregunta?

Hoy: Definir nuestro espacio de trabajo



```
getwd() # Directorio de trabajo actual  
setwd("ruta") # Establecer directorio de trabajo
```



Estructura de una base de datos



Las filas representan las **observaciones** - Las columnas representan las **variables**

```
id <- c(1, 2, 3, 4)
edad <- c(23, 45, 67, 89)
sexo <- c(1, 0, 1, 0)
peso <- c(80, 60, 70, 50)
altura <- c(180, 160, 200, 140)

data <- as.data.frame(cbind(id, edad, sexo,
                             peso, altura))

data
```

```
##   id edad sexo peso altura
## 1  1  23    1  80   180
## 2  2  45    0  60   160
## 3  3  67    1  70   200
## 4  4  89    0  50   140
```

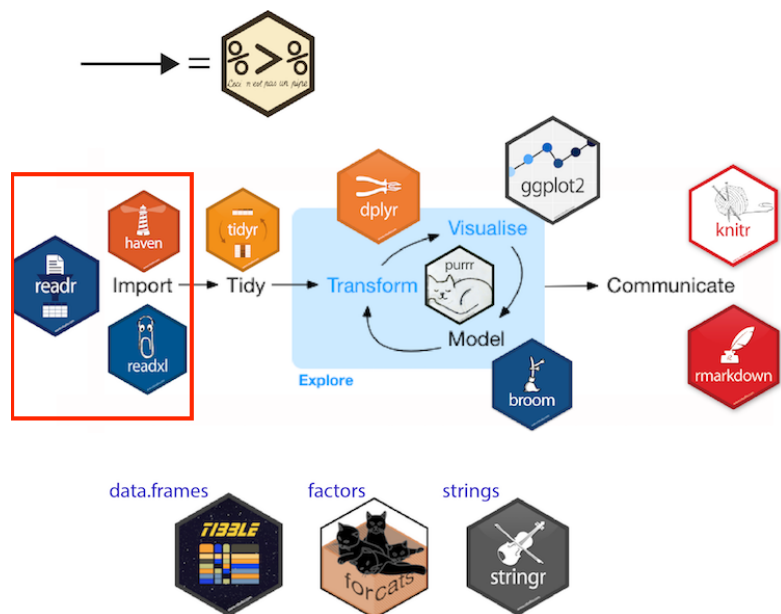
| Dim | Homogeneous | Heterogeneous |
|-----|---------------|---------------|
| 1d | Atomic Vector | Matrix |
| 2d | Matrix | Data - Frame |
| nd | Array | |

Mayor detalle pueden revisar el siguiente [enlace](#)

Importar una base de datos



```
install.packages(c("haven", "readxl", "readr"))  
library(haven)  
library(readxl)  
library(readr)
```



Mayor detalle del paquete [readr](#) aquí

Mayor detalle del paquete [haven](#) aquí

Explorar nuestra base de datos



```
dim(data)  # Observaciones y variables
```

```
## [1] 4 5
```

```
names(data) # Nombre de nuestras variables
```

```
## [1] "id"      "edad"    "sexo"    "peso"    "altura"
```

```
str(data)  # Visor de nuestras variables
```

```
## 'data.frame':    4 obs. of  5 variables:
## $ id      : num  1 2 3 4
## $ edad    : num  23 45 67 89
## $ sexo    : num  1 0 1 0
## $ peso    : num  80 60 70 50
## $ altura  : num  180 160 200 140
```

Explorar nuestra base de datos



```
head(data) # Primeras 6 observaciones
```

```
##   id edad sexo peso altura
## 1   1  23   1   80   180
## 2   2  45   0   60   160
## 3   3  67   1   70   200
## 4   4  89   0   50   140
```

```
tail(data) # Últimas 6 observaciones
```

```
##   id edad sexo peso altura
## 1   1  23   1   80   180
## 2   2  45   0   60   160
## 3   3  67   1   70   200
## 4   4  89   0   50   140
```

Explorar nuestra base de datos



```
library(skimr)
skim(data)
```

| skim_type | skim_variable | n_missing | numeric.mean | numeric.sd | numeric.p0 | numeric.p25 | numeric.p50 | numeric.p75 | numeric.p100 | numeric.hist |
|-----------|---------------|-----------|--------------|------------|------------|-------------|-------------|-------------|--------------|--------------|
| numeric | id | 0 | 2.5 | 1.2909944 | 1 | 1.75 | 2.5 | 3.25 | 4 | |
| numeric | edad | 0 | 56.0 | 28.4018779 | 23 | 39.50 | 56.0 | 72.50 | 89 | |
| numeric | sexo | 0 | 0.5 | 0.5773503 | 0 | 0.00 | 0.5 | 1.00 | 1 | |
| numeric | peso | 0 | 65.0 | 12.9099445 | 50 | 57.50 | 65.0 | 72.50 | 80 | |
| numeric | altura | 0 | 170.0 | 25.8198890 | 140 | 155.00 | 170.0 | 185.00 | 200 | |

¿Qué podemos ver de nuestras variables?

Validación de una base de datos



- Identificar observaciones duplicadas
- Revisar valores fuera de rango
- Explorar codificación de los casos perdidos
- Chequear efectividad de los filtros (si es que los hay)

Exportar una base de datos



Podemos guardar nuestra base como texto plano:

```
write.table(data, file="data.txt", sep="\t")
```

Podemos guardar en un archivo separado por comas (csv):

```
# Separado por comas
write.csv(data, file="data.csv", row.names = FALSE)

# Otras opciones de separación
write.csv(data, file="data.csv", row.names = FALSE, sep = "-")
```

¿Qué ventaja tienen estos formatos?

Podemos guardar en un excel tradicional:

```
library(openxlsx)
write.xlsx(data, file="data.xlsx")
```

Exportar una base de datos



Podemos guardar en el formato de [SPSS](#):

```
library(foreign)
write.foreign(data, datafile="data.txt", codefile="data.sps" package="SPSS")
```

Podemos guardar en el formato de [SAS](#):

```
write.foreign(data, datafile="data.txt", codefile="data.sas" package="SAS")
```

Podemos guardar en el formato de [Stata](#):

```
write.dta(data, datafile="data.dta")
```

Podemos guardar en el formato de [R](#):

```
save(rm, file="data.Rdata")
```

Próximo taller



- Encadenar operaciones: `pipe %>%`
- Filtrar y ordenar datos: `filter()` & `arrange`
- Seleccionar columnas: `select()`
- Generar nuevas variables: `mutate` & operador `$`
- Recodificar y etiquetar: `indexar code` & `dplyr`)

Próximo taller



Data Transformation with dplyr : : CHEAT SHEET

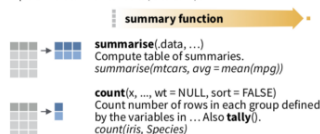


dplyr functions work with pipes and expect tidy data. In tidy data:



Summarise Cases

These apply **summary functions** to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).



VARIATIONS

summarise_all() - Apply funs to every column.
summarise_at() - Apply funs to specific columns.
summarise_if() - Apply funs to all cols of one type.

Group Cases

Use **group_by()** to create a "grouped" copy of a table. dplyr functions will manipulate each "group" separately and then combine the results.



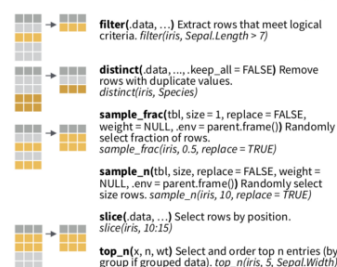
group_by(data, ..., add = FALSE)
Returns copy of table grouped by ...
`g_iris <- group_by(iris, Species)`

ungroup(x, ...)
Returns ungrouped copy of table.
`ungroup(g_iris)`

Manipulate Cases

EXTRACT CASES

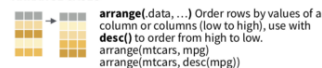
Row functions return a subset of rows as a new table.



Logical and boolean operators to use with filter()

`<` `<=` `is.na()` `%in%` `|` `xor()`
`>` `>=` `!is.na()` `!` `&`
See ?base::logic and ?Comparison for help.

ARRANGE CASES



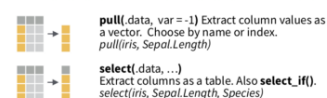
ADD CASES

add_row(data, ..., before = NULL, after = NULL)
Add one or more rows to a table.
`add_row(faithful, eruptions = 1, waiting = 1)`

Manipulate Variables

EXTRACT VARIABLES

Column functions return a set of columns as a new vector or table.

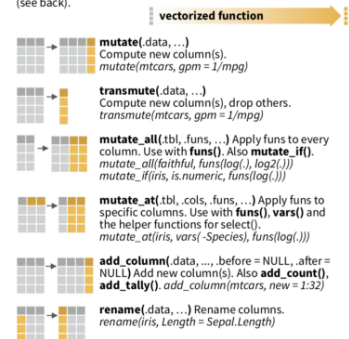


Use these helpers with **select()**, e.g. `select(iris, starts_with("Sepal"))`

contains(match) **num_range(prefix, range)** **ends_with(match)** **one_of(...)** **matches(match)** **starts_with(match)**
e.g. `mpg:cyl`
e.g. `Species`

MAKE NEW VARIABLES

These apply **vectorized functions** to columns. Vectorized funs take vectors as input and return vectors of the same length as output (see back).



RStudio is a trademark of RStudio, Inc. • CC BY SA RStudio • info@rstudio.com • 844-448-1212 • rstudio.com • Learn more with browseVignettes(package = "dplyr", "tidbale") • dplyr 0.7.0 • tidbale 1.2.0 • Updated: 2017-03

Pueden adelantar revisando el siguiente [enlace](#)

Donación



Se donará el 25% de lo recaudado en el taller: \$350.000.

Se está gestionando **cajas de alimentos** para personas con problemas económicos en La comuna de la Pintana, sector el castillo.

¡Se les avisará cuando se entreguen estas cajas! (en las **próximas 2 semanas**)

Sugerencias: ✉ **Buzón anónimo**

¡Gracias!

Presentación generada en:



✉ jdconejeros@uc.cl  [joseconejerosp](#)

 [JDConejeros](#)  [Jose_Conejeros](#)