

Café com estatística e R

Treinamento 2 - Importação e manipulação de dados no R e estatística descritiva

Marcelo Teixeira Paiva

2025-10-08

Abstract

Relatório do segundo treinamento onde foi apresentado como importar dados e manipulá-los no R, bem como as principais estatísticas descritivas univariadas e multivariadas.

Índice

1	Importação de dados no R	4
1.1	Pacotes necessários	4
1.2	Leitura de datasets externos ao R	8
1.2.1	Importando dados do Excel	8
1.2.2	Importando dados do Stata	11
1.2.3	Verificação e diagnóstico dos dados importados	13

Lista de Figuras

Lista de Tabelas

```
# Pacotes
library(tidyverse)
library(gridExtra)
library(plotly)
library(gt)
library(tidyverse)
library(kableExtra)

# Tema personalizado para gráficos
tema_didatico <- theme_minimal() +
  theme(
    plot.title = element_text(face = "bold", size = 14),
    plot.subtitle = element_text(face = "italic", size = 11),
    axis.title = element_text(size = 12),
    legend.position = "top",
    panel.grid.minor = element_blank()
  )

cores <- c("#FF6B6B", "#4ECDC4", "#45B7D1", "#96CEB4", "#FFEAA7")
```

Chapter 1

Importação de dados no R

1.1 Pacotes necessários

Pacotes (**package**) são coleções de funções, dados e documentação que estendem as capacidades do R base (aquele que você recebe na instalação padrão). São como “caixas de ferramentas” especializadas que você adiciona ao R para realizar tarefas específicas, então você tem pacotes para elaboração de gráficos, para certos tipos de análises, para manipulação de dados, para leitura (importação) de dados. Em <https://cran.r-project.org/web/views/> há uma “breve” lista de pacotes conforme a sua finalidade.

```
# funções no R base
length(ls("package:base"))
```

```
[1] 1270
```

```
# funções especializadas no pacote dplyr
length(ls("package:dplyr"))
```

```
[1] 297
```

```
# um pacote possui um conjunto de arquivos associados
system.file(package = "ggplot2") %>% list.files()
```

```
[1] "CITATION"      "data"          "DESCRIPTION"   "doc"           "help"
[6] "html"          "INDEX"         "LICENSE"       "Meta"          "NAMESPACE"
[11] "NEWS.md"       "R"
```

Por padrão, ao iniciar uma sessão no R, serão carregados os pacotes e funções associados ao R base. Os demais devem ser instalados primeiramente, e depois carregados na seção para serem usados.

```
# pacotes carregados no seu ambiente
search()
```

```
[1] ".GlobalEnv"      "package:kableExtra" "package:gt"
[4] "package:plotly"  "package:gridExtra"  "package:lubridate"
[7] "package:forcats" "package:stringr"    "package:dplyr"
[10] "package:purrr"   "package:readr"      "package:tidyr"
[13] "package:tibble"  "package:ggplot2"    "package:tidyverse"
[16] "package:stats"   "package:graphics"   "package:grDevices"
[19] "package:utils"   "package:datasets"   "package:methods"
```

```
[22] "Autoloads"          "package:base"

# pacotes instalados
instalados <- installed.packages()[, "Package"]
instalados[1:4]

      abind      ARTool      askpass      backports
"abind"    "ARTool"    "askpass" "backports"

length(instalados)

[1] 330

# verificando se um pacote já está instalado
sum(installed.packages()[, "Package"] == 'dplyr')

[1] 1

any(installed.packages()[, "Package"] == 'dplyr')

[1] TRUE

"ggplot2" %in% rownames(installed.packages())

[1] TRUE
```

A instalação de pacotes no R é feita usando a função `install.packages` ou `devtools::install_github` para pacotes que estão no github e não em um repositório de pacotes.

```
# pelo repositório oficial (na web)
install.packages("ggplot2")
install.packages(c("dplyr", "tidyr", "readr")) # instalando vários pacotes de uma vez

# Instalar o pacote e todas dependências relacionadas a ele
install.packages("ggplot2", dependencies = TRUE)

# instalar de um arquivo local
install.packages("caminho/para/pacote.tar.gz", repos = NULL, type = "source")

# Instalar pacote mantido no GitHub
install.packages("devtools")
devtools::install_github("tidyverse/ggplot2")

# Usar outros repositórios para instalação
install.packages("ggplot2", repos = "https://cloud.r-project.org/")
```

Para carregar um pacote em uma sessão usamos `library()` ou `require()`. A diferença entre os dois é que, na ausência do pacote que você pretende carregar, `library` gera um erro, enquanto o `require` retorna um valor `FALSE` invisível, o qual pode ser usado, por exemplo, para criar uma lógica em seu script para instalar o pacote caso o mesmo não possa ser carregado ou, então, para gerar uma mensagem no terminal indicando essa ausência do pacote.

```
library(ggplot2)

# Não exibir mensagens de carregamento do pacote
suppressPackageStartupMessages(library(ggplot2))
```

```
# criando uma lógica simples com require para instalar pacotes que
# não possam ser carregados
if(!require(ggplot2)) {
  install.packages("ggplot2")
  require(ggplot2)
}

# usando uma função do pacote sem o carregar (namespace qualification)
head(dplyr::filter(mtcars, mpg > 20), 2)
```

```
      mpg cyl disp  hp drat   wt  qsec vs am gear carb
Mazda RX4     21   6  160 110  3.9 2.620 16.46  0  1   4    4
Mazda RX4 Wag 21   6  160 110  3.9 2.875 17.02  0  1   4    4
```

```
# carregando vários pacotes de uma lista de nomes
pacotes <- c("ggplot2", "dplyr", "tidyr")
x <- lapply(pacotes, library, character.only = TRUE, quietly = TRUE)
```

Além dessas funções para instalação e carregamento de pacotes, também outras funções que devem ser conhecidas na rotina são as de atualização (`update.packages()`) e remoção (`remove.packages()`) de pacotes, descrição (`packageDescription()`), versão (`packageVersion()`) e forma recomendada pelo seus autores de citação (`citation()`) quando usada em uma publicação.

```
# Atualização de pacotes
update.packages() # todos
update.packages(ask = FALSE) # todos, mas exige confirmação

# apagar um pacote
remove.packages("nome_pacote")

# descrição e versão
packageDescription("ggplot2")
```

```
Package: ggplot2
Title: Create Elegant Data Visualisations Using the Grammar of Graphics
Version: 4.0.0
Authors@R: c( person("Hadley", "Wickham", , "hadley@posit.co", role =
  "aut", comment = c(ORCID = "0000-0003-4757-117X")),
  person("Winston", "Chang", role = "aut", comment = c(ORCID =
  "0000-0002-1576-2126")), person("Lionel", "Henry", role =
  "aut"), person("Thomas Lin", "Pedersen", ,
  "thomas.pedersen@posit.co", role = c("aut", "cre"), comment =
  c(ORCID = "0000-0002-5147-4711")), person("Kohske",
  "Takahashi", role = "aut"), person("Claus", "Wilke", role =
  "aut", comment = c(ORCID = "0000-0002-7470-9261")),
  person("Kara", "Woo", role = "aut", comment = c(ORCID =
  "0000-0002-5125-4188")), person("Hiroaki", "Yutani", role =
  "aut", comment = c(ORCID = "0000-0002-3385-7233")),
  person("Dewey", "Dunnington", role = "aut", comment = c(ORCID =
  "0000-0002-9415-4582")), person("Teun", "van den Brand", role =
  "aut", comment = c(ORCID = "0000-0002-9335-7468")),
  person("Posit, PBC", role = c("cph", "fnd"), comment = c(ROR =
```



```

"03wc8by49")) )
Description: A system for 'declaratively' creating graphics, based on
"The Grammar of Graphics". You provide the data, tell 'ggplot2'
how to map variables to aesthetics, what graphical primitives
to use, and it takes care of the details.
License: MIT + file LICENSE
URL: https://ggplot2.tidyverse.org,
https://github.com/tidyverse/ggplot2
BugReports: https://github.com/tidyverse/ggplot2/issues
Depends: R (>= 4.1)
Imports: cli, grDevices, grid, gtable (>= 0.3.6), isoband, lifecycle (>
1.0.1), rlang (>= 1.1.0), S7, scales (>= 1.4.0), stats, vctrs
(>= 0.6.0), withr (>= 2.5.0)
Suggests: broom, covr, dplyr, ggplot2movies, hexbin, Hmisc, knitr,
mapproj, maps, MASS, mgcv, multcomp, munsell, nlme, profvis,
quantreg, ragg (>= 1.2.6), RColorBrewer, rmarkdown, roxygen2,
rpart, sf (>= 0.7-3), svglite (>= 2.1.2), testthat (>= 3.1.5),
tibble, vdiffr (>= 1.0.6), xml2
Enhances: sp
VignetteBuilder: knitr
Config/Needs/website: ggtext, tidyr, forcats, tidyverse/tidytemplate
Config/testthat/edition: 3
Config/usethis/last-upkeep: 2025-04-23
Encoding: UTF-8
LazyData: true
RoxygenNote: 7.3.2
Collate: 'ggproto.R' 'ggplot-global.R' 'aaa-.R'
'aes-colour-fill-alpha.R' .....
NeedsCompilation: no
Packaged: 2025-08-19 08:21:45 UTC; thomas
Author: Hadley Wickham [aut] (ORCID:
<https://orcid.org/0000-0003-4757-117X>), Winston Chang [aut]
(ORCID: <https://orcid.org/0000-0002-1576-2126>), Lionel Henry
[aut], Thomas Lin Pedersen [aut, cre] (ORCID:
<https://orcid.org/0000-0002-5147-4711>), Kohske Takahashi
[aut], Claus Wilke [aut] (ORCID:
<https://orcid.org/0000-0002-7470-9261>), Kara Woo [aut]
(ORCID: <https://orcid.org/0000-0002-5125-4188>), Hiroaki
Yutani [aut] (ORCID: <https://orcid.org/0000-0002-3385-7233>),
Dewey Dunnington [aut] (ORCID:
<https://orcid.org/0000-0002-9415-4582>), Teun van den Brand
[aut] (ORCID: <https://orcid.org/0000-0002-9335-7468>), Posit,
PBC [cph, fnd] (ROR: <https://ror.org/03wc8by49>)
Maintainer: Thomas Lin Pedersen <thomas.pedersen@posit.co>
Repository: CRAN
Date/Publication: 2025-09-11 07:10:02 UTC
Built: R 4.3.3; ; 2025-10-07 18:16:38 UTC; unix

-- File: /home/marcelo/R/x86_64-pc-linux-gnu-library/4.3/ggplot2/Meta/package.rds

```

```
packageVersion("ggplot2")
```

```
[1] '4.0.0'
```

```
# forma de citação  
citation("ggplot2")
```

To cite ggplot2 in publications, please use

H. Wickham. ggplot2: Elegant Graphics for Data Analysis.
Springer-Verlag New York, 2016.

A BibTeX entry for LaTeX users is

```
@Book{,  
  author = {Hadley Wickham},  
  title = {ggplot2: Elegant Graphics for Data Analysis},  
  publisher = {Springer-Verlag New York},  
  year = {2016},  
  isbn = {978-3-319-24277-4},  
  url = {https://ggplot2.tidyverse.org},  
}
```

Algo a se ter em mente é que nada impede de vários pacotes terem o mesmo nome para funções com finalidades diferentes. Nesse caso, ao carregar esses pacotes, o último a ser carregado irá mascarar o nome da anterior no seu ambiente. Assim, para evitar conflitos, ou o uso da função errada, recomenda-se usar a função seguindo o padrão `nome_do_pacote::nome_da_função`.

1.2 Leitura de datasets externos ao R

A importação de dados é o primeiro passo em qualquer análise. O R oferece múltiplos pacotes especializados para diferentes formatos de arquivos, mas iremos focar nos pacotes de leitura dos arquivos provenientes dos softwares Excel, SAS, Stata e SPSS. Para isso, utilizaremos os pacotes `readxl` e `haven`.

```
# mini rotina para instalar um pacote se ainda não estiver instalado  
instala_se_nao_existe <- function(nome_do_pacote){  
  if(nome_do_pacote %in% rownames(installed.packages())) return()  
  install.packages(nome_do_pacote)  
}  
lapply(c("readxl", "haven"), instala_se_nao_existe)  
  
# Carregar pacotes  
library(readxl) # Excel  
library(haven)  # SAS, SPSS, STATA
```

1.2.1 Importando dados do Excel

Para leitura de arquivos do Excel nos formatos `.xls` e `.xlsx` usaremos o pacote `readxl`, o qual faz parte do conjunto de pacotes do `tidyverse`. Dele podemos usar as funções `read_excel()`, `read_xls()` ou `read_xlsx()`, os quais recebem argumentos semelhantes, com a diferença que os dois últimos são específicos ao formato do arquivo.

O primeiro e mais importante argumento a ser fornecido para essa função é o **path**, o local onde o arquivo se encontra no seu computador. Esse caminho pode ser absoluto (desde a raiz, normalmente / no linux ou C: no windows, até o local) ou relativo ao diretório de trabalho (que pode ser verificado usando a função `getwd()`).

Como os arquivos do Excel aceitam múltiplas planilhas (em diferentes abas), o argumento de **sheet** do `read_excel()` permite escolher qual aba se pretende carregar. Caso seja necessário verificar primeiro o nome das abas disponíveis no arquivo, use `excel_sheets(path)`.

Outro problema comum em arquivos do Excel são planilhas que não iniciam na linha 1 ou que apresentam um conjunto de colunas que não pretendemos usar (sem conteúdo ou preenchido com informações que não fazem parte do dataset). Para contornar esses obstáculos, podemos usar o argumento **skip** com o número de linhas iniciais que não devem ser lidas, ou usar o **range** com um **character** indicando a primeira e última células que delimitam seus dados (por exemplo, `range = "B2:D20"` indica que devem ser lidas as colunas B, C e D, das linhas 2 até a 20).

Por padrão, essas funções buscam adivinhar o tipo de dados presente em cada coluna da planilha, mas é possível declarar o tipo usando o argumento `col_types` com um vetor com comprimento igual ao número de colunas que irá importar. Esse vetor deve, para cada coluna, usar uma das opções:

- “skip”: remove a coluna do dataset
- “guess”: deixa para a função escolher o tipo
- “logical”: booleano
- “numeric”: numérico
- “date”: data
- “text”: character
- “list”: lista

Também por padrão, a primeira linha é usada para obter os nomes de cada coluna. Se você não possui nomes das colunas na sua planilha use `col_names = FALSE` na função ou passe um vetor dos nomes das colunas para o argumento `col_names`.

Um aspecto importante de qualquer conjunto de dados é saber como foram codificados os dados ausentes. O argumento **na** permite passar um vetor de **character** com os códigos usados na planilha para declarar um dado ausente, o qual será convertido para NA no R.

```
excel_sheets("../datasets/excel/ap2.xlsx")
```

```
[1] "Data"
```

```
dados_excel <- read_excel("../datasets/excel/ap2.xlsx")
head(dados_excel)
```

```
# A tibble: 6 x 21
  farm_id batch_id litt_id pig_id parity vacc_mp seas_fin age_t w_age_t age_t6
  <dbl>    <dbl>    <dbl> <dbl>   <dbl>   <dbl>    <dbl> <dbl>   <dbl> <dbl>
1      1      1      1      1      1      8      1      1     70    33.8    116
2      1      1      1      1      2      8      1      1     70    32.9    116
3      1      1      1      1      3      8      1      1     70    29.4    116
4      1      1      1      2      4      8      1      1     60    19.8    106
5      1      1      1      2      5      8      1      1     60    20.4    106
6      1      1      1      2      6      8      1      1     60    20.3    106
# i 11 more variables: w_age_t6 <dbl>, dwg_fin <dbl>, ap2_t <dbl>, mp_t <dbl>,
#   infl_t <dbl>, prrs_t <dbl>, ap2_t6 <dbl>, mp_t6 <dbl>, infl_t6 <dbl>,
#   prrs_t6 <dbl>, ap2_sc <dbl>
```

```
# definir a planilha por nome ou índice
dados_pela_aba <- read_excel("../datasets/excel/ap2.xlsx", sheet = "Data")
dados_pela_aba <- read_excel("../datasets/excel/ap2.xlsx", sheet = 1)

# carregar somente um intervalo de células, em que a linha 1 não é header
dados_pelo_range <- read_excel(
  "../datasets/excel/ap2.xlsx",
  range = "A2:B100",
  sheet = "Data",
  col_names = FALSE
)
```

New names:

```
* `` -> `...1`
* `` -> `...2`
```

```
head(dados_pelo_range)
```

```
# A tibble: 6 x 2
  ...1 ...2
  <dbl> <dbl>
1     1     1
2     1     1
3     1     1
4     1     1
5     1     1
6     1     1
```

mesmo exemplo, mas definindo os nomes das colunas

```
dados_pelo_range <- read_excel(
  "../datasets/excel/ap2.xlsx",
  range = "A2:B100",
  sheet = "Data",
  col_names = c('fazenda', 'lote')
)
head(dados_pelo_range)
```

```
# A tibble: 6 x 2
  fazenda lote
  <dbl> <dbl>
1     1     1
2     1     1
3     1     1
4     1     1
5     1     1
6     1     1
```

declarando os tipos de colunas

```
dados_tipos <- read_excel(
  "../datasets/excel/ap2.xlsx",
  col_types = c("text", "numeric", rep("text", 19))
)
```

```
head(dados_tipos)
```

```
# A tibble: 6 x 21
  farm_id batch_id litt_id pig_id parity vacc_mp seas_fin age_t w_age_t age_t6
  <chr>      <dbl> <chr>   <chr>   <chr>   <chr>   <chr>   <chr> <chr>   <chr>
1 1          1 1       1       8       1       1       70  33.8   116
2 1          1 1       2       8       1       1       70  32.9   116
3 1          1 1       3       8       1       1       70  29.4   116
4 1          1 2       4       8       1       1       60  19.8   106
5 1          1 2       5       8       1       1       60  20.4   106
6 1          1 2       6       8       1       1       60  20.3   106
# i 11 more variables: w_age_t6 <chr>, dwg_fin <chr>, ap2_t <chr>, mp_t <chr>,
#   infl_t <chr>, prrs_t <chr>, ap2_t6 <chr>, mp_t6 <chr>, infl_t6 <chr>,
#   prrs_t6 <chr>, ap2_sc <chr>
```

```
# definir os códigos usados na planilha para dados ausentes
```

```
dados_na <- read_excel(
  "../datasets/excel/ap2.xlsx",
  na = c("", "NA", "N/A", "-")
)
```

1.2.2 Importando dados do Stata

Para leitura de arquivos do Stata no formato .dta usaremos o pacote **heaven**, o qual possui funções para leitura de arquivos do Stata, SPSS e SAS. Nesse treinamento vamos focar na função `read_dta()` para leitura dos arquivos do Stata (superiores a versão 13.0).

Assim como no `read_excel()`, o primeiro argumento de `read_dta()` deve ser a localização do arquivo. Além disso, a função aceita como argumentos `encoding`, a codificação de caracteres usada, `skip` para remover um certo número de linhas, `col_select` para definir quais colunas serão seleccionadas e `n_max` para declarar o número máximo de linhas que devem ser importadas.

Um diferença importante entre arquivos do Excel e do Stata é que no segundo o dataset e as suas variáveis podem conter metadados (“notes” e “labels”) com informações sobre esses dados. Essas informações podem ser acessadas na função `attr()`.

```
dados_stata <- read_dta("../datasets/stata/ap2.dta")
head(dados_stata)
```

```
# A tibble: 6 x 21
  farm_id batch_id litt_id pig_id parity vacc_mp seas_fin age_t w_age_t age_t6
  <dbl>      <dbl>   <dbl>   <dbl>   <dbl> <dbl+lbl> <dbl+lbl> <dbl>   <dbl>   <dbl>
1      1          1       1       1       8 1 [vac]  1 [wint~  70  33.8   116
2      1          1       1       2       8 1 [vac]  1 [wint~  70  32.9   116
3      1          1       1       3       8 1 [vac]  1 [wint~  70  29.4   116
4      1          1       2       4       8 1 [vac]  1 [wint~  60  19.8   106
5      1          1       2       5       8 1 [vac]  1 [wint~  60  20.4   106
6      1          1       2       6       8 1 [vac]  1 [wint~  60  20.3   106
# i 11 more variables: w_age_t6 <dbl>, dwg_fin <dbl>, ap2_t <dbl+lbl>,
#   mp_t <dbl+lbl>, infl_t <dbl+lbl>, prrs_t <dbl+lbl>, ap2_t6 <dbl+lbl>,
#   mp_t6 <dbl+lbl>, infl_t6 <dbl+lbl>, prrs_t6 <dbl+lbl>, ap2_sc <dbl+lbl>
```

```

dados_stata_com_encoding <- read_dta(
  "../datasets/stata/ap2.dta",
  encoding = "UTF-8"
)

# transformar colunas labelled em factor
dados_stata_como_factor <- read_dta(
  "../datasets/stata/ap2.dta",
  encoding = "UTF-8"
) |> as_factor()
head(dados_stata_como_factor)

# A tibble: 6 x 21
  farm_id batch_id litt_id pig_id parity vacc_mp seas_fin age_t w_age_t age_t6
    <dbl>   <dbl>   <dbl>  <dbl>  <dbl> <fct>   <fct>   <dbl>   <dbl>   <dbl>
1       1       1       1      1      1     8 vac   winter    70    33.8    116
2       1       1       1      2      2     8 vac   winter    70    32.9    116
3       1       1       1      3      3     8 vac   winter    70    29.4    116
4       1       1       2      4      4     8 vac   winter    60    19.8    106
5       1       1       2      5      5     8 vac   winter    60    20.4    106
6       1       1       2      6      6     8 vac   winter    60    20.3    106
# i 11 more variables: w_age_t6 <dbl>, dwg_fin <dbl>, ap2_t <fct>, mp_t <fct>,
#   infl_t <fct>, prrs_t <fct>, ap2_t6 <fct>, mp_t6 <fct>, infl_t6 <fct>,
#   prrs_t6 <fct>, ap2_sc <fct>

# Notas do Stata
attr(dados_stata, "notes")

[1] "5 Aug 2002 16:24 data provided by Dr. Haakan Vigre, Denmark"
[2] "1"

# Labels das variáveis
labels <- sapply(dados_stata, function(x) attr(x, "label"))
kable(
  tibble(var=names(labels), metadata=labels),
  col.names = c("Variável", "Label")
)

```

Variável	Label
farm_id	farm identification
batch_id	batch identifiaction number
litt_id	litter identification number
pig_id	pig identification
parity	the farrowing no. of the sow
vacc_mp	the batch vaccinated against M.hyop yes=1
seas_fin	prod. season in finishing unit: winther=1
age_t	pig-age transfer from weaning to finishing unit
w_age_t	weight in kg. at age_t
age_t6	age_tra plus approx. 6 weeks
w_age_t6	weight in kg. at age_t6
dwg_fin	dwg in g. between age_t and age_t6

Variável	Label
ap2_t	serological reac. against A.pleuropneumoniae serotype 2 at age_t
mp_t	serological reac. against M.hyopneumoniae at age_t
infl_t	serological reac. against Influenza virus at age_t
prrs_t	serological reac. against PRRS virus at age_t
ap2_t6	serological reac. against A.pleuropneumoniae serotype 2 at age_t6
mp_t6	serological reac. against M.hyopneumoniae at age_t6
infl_t6	serological reac. against Influenza virus at age_t6
prrs_t6	serological reac. against PRRS virus at age_t6
ap2_sc	seroconversion to ap2 during the finishing period

1.2.3 Verificação e diagnóstico dos dados importados

Uma vez carregados os dados, é importante avaliar a estrutura desse conjunto de dados importado. Para uma exploração inicial, será interessante avaliar, no mínimo, as dimensões desses dados (número de observações e variáveis), quais os tipos das variáveis no R, resumos estatísticos simples, quantidade de valores ausentes por variável.

```
verificar_dados <- function(dados) {
  cat("Dimensões:", dim(dados), "\n")
  cat("Tipos de variáveis:\n")
  print(sapply(dados, class))
  cat("\nPrimeiras linhas:\n")
  print(head(dados, 3))
  cat("\nResumo estatístico:\n")
  print(summary(dados))
  cat("\nValores missing por coluna:\n")
  print(colSums(is.na(dados)))
  cat("\nStructura dos dados:\n")
  str(dados)
}

# Aplicar a qualquer dataset importado
verificar_dados(dados_stata)
```

```
Dimensões: 1114 21
Tipos de variáveis:
$farm_id
[1] "numeric"

$batch_id
[1] "numeric"

$litt_id
[1] "numeric"

$pig_id
[1] "numeric"

$parity
```

```

[1] "numeric"

$vaccc_mp
[1] "haven_labelled" "vctrs_vctr"      "double"

$seas_fin
[1] "haven_labelled" "vctrs_vctr"      "double"

$age_t
[1] "numeric"

$w_age_t
[1] "numeric"

$age_t6
[1] "numeric"

$w_age_t6
[1] "numeric"

$dwg_fin
[1] "numeric"

$ap2_t
[1] "haven_labelled" "vctrs_vctr"      "double"

$mp_t
[1] "haven_labelled" "vctrs_vctr"      "double"

$infl_t
[1] "haven_labelled" "vctrs_vctr"      "double"

$prrs_t
[1] "haven_labelled" "vctrs_vctr"      "double"

$ap2_t6
[1] "haven_labelled" "vctrs_vctr"      "double"

$mp_t6
[1] "haven_labelled" "vctrs_vctr"      "double"

$infl_t6
[1] "haven_labelled" "vctrs_vctr"      "double"

$prrs_t6
[1] "haven_labelled" "vctrs_vctr"      "double"

$ap2_sc
[1] "haven_labelled" "vctrs_vctr"      "double"

```


Primeiras linhas:

A tibble: 3 x 21

```

  farm_id batch_id litt_id pig_id parity vacc_mp seas_fin age_t w_age_t age_t6
    <dbl>   <dbl>   <dbl>  <dbl>  <dbl> <dbl+lbl> <dbl+lb> <dbl>   <dbl>   <dbl>
1       1       1       1      1      8 1 [vac]  1 [wint~   70   33.8   116
2       1       1       1      2      8 1 [vac]  1 [wint~   70   32.9   116
3       1       1       1      3      8 1 [vac]  1 [wint~   70   29.4   116
# i 11 more variables: w_age_t6 <dbl>, dwg_fin <dbl>, ap2_t <dbl+lbl>,
#   mp_t <dbl+lbl>, infl_t <dbl+lbl>, prrs_t <dbl+lbl>, ap2_t6 <dbl+lbl>,
#   mp_t6 <dbl+lbl>, infl_t6 <dbl+lbl>, prrs_t6 <dbl+lbl>, ap2_sc <dbl+lbl>

```

Resumo estatístico:

farm_id	batch_id	litt_id	pig_id
Min. :1.000	Min. : 1.0	Min. : 1.0	Min. : 1.0
1st Qu.:2.000	1st Qu.:11.0	1st Qu.:124.0	1st Qu.: 371.2
Median :3.000	Median :22.0	Median :256.5	Median : 766.5
Mean :3.273	Mean :21.6	Mean :252.2	Mean : 752.9
3rd Qu.:5.000	3rd Qu.:33.0	3rd Qu.:388.0	3rd Qu.:1156.8
Max. :6.000	Max. :41.0	Max. :491.0	Max. :1466.0

parity	vacc_mp	seas_fin	age_t
Min. : 1.000	Min. :0.0000	Min. :0.000	Min. :52.00
1st Qu.: 2.000	1st Qu.:0.0000	1st Qu.:0.000	1st Qu.:62.00
Median : 4.000	Median :1.0000	Median :0.000	Median :68.00
Mean : 3.769	Mean :0.6239	Mean :0.342	Mean :68.91
3rd Qu.: 5.000	3rd Qu.:1.0000	3rd Qu.:1.000	3rd Qu.:76.00
Max. :11.000	Max. :1.0000	Max. :1.000	Max. :83.00

w_age_t	age_t6	w_age_t6	dwg_fin
Min. :12.10	Min. :101.0	Min. : 30.00	Min. : 231.0
1st Qu.:22.30	1st Qu.:117.0	1st Qu.: 61.00	1st Qu.: 647.5
Median :26.80	Median :124.0	Median : 70.00	Median : 758.5
Mean :27.43	Mean :125.6	Mean : 69.97	Mean : 757.5
3rd Qu.:32.10	3rd Qu.:131.0	3rd Qu.: 78.00	3rd Qu.: 872.8
Max. :47.20	Max. :161.0	Max. :107.00	Max. :1188.0

ap2_t	mp_t	infl_t	prrs_t
Min. :0.0000	Min. :0.0000	Min. :0.0000	Min. :0.0000
1st Qu.:0.0000	1st Qu.:0.0000	1st Qu.:0.0000	1st Qu.:0.0000
Median :0.0000	Median :0.0000	Median :0.0000	Median :0.0000
Mean :0.1059	Mean :0.2738	Mean :0.1795	Mean :0.2127
3rd Qu.:0.0000	3rd Qu.:1.0000	3rd Qu.:0.0000	3rd Qu.:0.0000
Max. :1.0000	Max. :1.0000	Max. :1.0000	Max. :1.0000

ap2_t6	mp_t6	infl_t6	prrs_t6
Min. :0.0000	Min. :0.0000	Min. :0.0000	Min. :0.0000
1st Qu.:0.0000	1st Qu.:0.0000	1st Qu.:0.0000	1st Qu.:0.0000
Median :1.0000	Median :0.0000	Median :0.0000	Median :0.0000
Mean :0.5566	Mean :0.3293	Mean :0.4668	Mean :0.4102
3rd Qu.:1.0000	3rd Qu.:1.0000	3rd Qu.:1.0000	3rd Qu.:1.0000
Max. :1.0000	Max. :1.0000	Max. :1.0000	Max. :1.0000

NA's :33

```

ap2_sc
Min. :0.00
1st Qu.:0.00
Median :1.00
Mean :0.51
3rd Qu.:1.00
Max. :1.00
NA's :118

```

Valores missing por columna:

farm_id	batch_id	litt_id	pig_id	parity	vacc_mp	seas_fin	age_t
0	0	0	0	0	0	0	0
w_age_t	age_t6	w_age_t6	dwg_fin	ap2_t	mp_t	infl_t	prrs_t
0	0	0	0	0	0	0	0
ap2_t6	mp_t6	infl_t6	prrs_t6	ap2_sc			
0	33	0	0	118			

Structura dos dados:

```

tibble [1,114 x 21] (S3: tbl_df/tbl/data.frame)
 $ farm_id : num [1:1114] 1 1 1 1 1 1 1 1 1 1 ...
  .. attr(*, "label")= chr "farm identification"
  .. attr(*, "format.stata")= chr "%5.0f"
 $ batch_id: num [1:1114] 1 1 1 1 1 1 1 1 1 1 ...
  .. attr(*, "label")= chr "batch identifiacion number"
  .. attr(*, "format.stata")= chr "%5.0f"
 $ litt_id : num [1:1114] 1 1 1 2 2 2 3 3 3 4 ...
  .. attr(*, "label")= chr "litter identification number"
  .. attr(*, "format.stata")= chr "%5.0f"
 $ pig_id  : num [1:1114] 1 2 3 4 5 6 7 8 9 10 ...
  .. attr(*, "label")= chr "pig identification"
  .. attr(*, "format.stata")= chr "%5.0f"
 $ parity  : num [1:1114] 8 8 8 8 8 8 6 6 6 6 ...
  .. attr(*, "label")= chr "the farrowing no. of the sow"
  .. attr(*, "format.stata")= chr "%3.0f"
 $ vacc_mp : dbl+lbl [1:1114] 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ...
  ..@ label      : chr "the batch vaccinated against M.hyop yes=1"
  ..@ format.stata: chr "%8.0f"
  ..@ labels     : Named num [1:2] 0 1
  .. ..- attr(*, "names")= chr [1:2] "not vac." "vac"
 $ seas_fin: dbl+lbl [1:1114] 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ...
  ..@ label      : chr "prod. season in finishing unit: winther=1"
  ..@ format.stata: chr "%6.0f"
  ..@ labels     : Named num [1:2] 0 1
  .. ..- attr(*, "names")= chr [1:2] "summer" "winter"
 $ age_t   : num [1:1114] 70 70 70 60 60 60 67 67 67 61 ...
  .. attr(*, "label")= chr "pig-age transfer from weaning to finishing unit"
  .. attr(*, "format.stata")= chr "%5.0f"
 $ w_age_t : num [1:1114] 33.8 32.9 29.4 19.8 20.4 20.3 21 32.4 30.3 22.5 ...
  .. attr(*, "label")= chr "weight in kg. at age_t"
  .. attr(*, "format.stata")= chr "%5.1f"

```

```

$ age_t6 : num [1:1114] 116 116 116 106 106 106 113 113 113 107 ...
..- attr(*, "label")= chr "age_tra plus approx. 6 weeks"
..- attr(*, "format.stata")= chr "%5.0f"
$ w_age_t6: num [1:1114] 80 80 81 54 64 63 59 79 72 66 ...
..- attr(*, "label")= chr "weight in kg. at age_t6"
..- attr(*, "format.stata")= chr "%5.1f"
$ dwg_fin : num [1:1114] 1004 1024 1122 743 948 ...
..- attr(*, "label")= chr "dwg in g. between age_t and age_t6"
..- attr(*, "format.stata")= chr "%5.0f"
$ ap2_t : dbl+lbl [1:1114] 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,...
..@ label : chr "serological reac. against A.pleuropneumoniae serotype 2 at age_t"
..@ format.stata: chr "%3.0f"
..@ labels : Named num [1:2] 0 1
.. ..- attr(*, "names")= chr [1:2] "neg" "pos"
$ mp_t : dbl+lbl [1:1114] 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 0, 0,...
..@ label : chr "serological reac. against M.hypopneumoniae at age_t"
..@ format.stata: chr "%3.0f"
..@ labels : Named num [1:2] 0 1
.. ..- attr(*, "names")= chr [1:2] "neg" "pos"
$ infl_t : dbl+lbl [1:1114] 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,...
..@ label : chr "serological reac. against Influenza virus at age_t"
..@ format.stata: chr "%3.0f"
..@ labels : Named num [1:2] 0 1
.. ..- attr(*, "names")= chr [1:2] "neg" "pos"
$ prrs_t : dbl+lbl [1:1114] 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1,...
..@ label : chr "serological reac. against PRRS virus at age_t"
..@ format.stata: chr "%3.0f"
..@ labels : Named num [1:2] 0 1
.. ..- attr(*, "names")= chr [1:2] "neg" "pos"
$ ap2_t6 : dbl+lbl [1:1114] 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,...
..@ label : chr "serological reac. against A.pleuropneumoniae serotype 2 at age_t6"
..@ format.stata: chr "%3.0f"
..@ labels : Named num [1:2] 0 1
.. ..- attr(*, "names")= chr [1:2] "neg" "pos"
$ mp_t6 : dbl+lbl [1:1114] 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1, 0, 1, 0, 0,...
..@ label : chr "serological reac. against M.hypopneumoniae at age_t6"
..@ format.stata: chr "%3.0f"
..@ labels : Named num [1:2] 0 1
.. ..- attr(*, "names")= chr [1:2] "neg" "pos"
$ infl_t6 : dbl+lbl [1:1114] 0, 1, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 1, 1, 1, 1,...
..@ label : chr "serological reac. against Influenza virus at age_t6"
..@ format.stata: chr "%3.0f"
..@ labels : Named num [1:2] 0 1
.. ..- attr(*, "names")= chr [1:2] "neg" "pos"
$ prrs_t6 : dbl+lbl [1:1114] 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,...
..@ label : chr "serological reac. against PRRS virus at age_t6"
..@ format.stata: chr "%3.0f"
..@ labels : Named num [1:2] 0 1
.. ..- attr(*, "names")= chr [1:2] "neg" "pos"
$ ap2_sc : dbl+lbl [1:1114] 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,...
..@ label : chr "seroconversion to ap2 during the finishing period"

```

```

..@ format.stata: chr "%3.0f"
..@ labels      : Named num [1:2] 0 1
.. ..- attr(*, "names")= chr [1:2] "neg" "pos"
- attr(*, "notes")= chr [1:2] "5 Aug 2002 16:24 data provided by Dr. Haakan Vigre, Denmark" "1"

```

Por fim, em grandes datasets é comum que os dados sejam registrados em múltiplos arquivos (principalmente no Excel, por causa do limite de linhas). Nesse caso, para não ser necessário carregar cada um desses arquivos e depois construir um data.frame que uni todos, podemos usar recursos de programação funcional do pacote `purrr` para importar diretamente todos os arquivos em um único data.frame.

```

library(purrr)
# obter uma lista dos arquivos que serão importados
arquivos <- list.files("datasets/csv", pattern = "\\..csv$", full.names = TRUE)
# mapear todos os arquivos para um unico data.frame
dados <- map_df(arquivos, read_csv2)

```

Quadro Resumo das funções que podem ser usadas na importação de arquivos externos ao R

Formato	Pacote	Função
CSV	readr	read_delim(), read_csv(), read_csv2()
Excel	readxl	read_excel(), read_xls(), read_xlsx()
SAS	haven	read_sas()
SPSS	haven	read_sav()
Stata	haven	read_stata(), read_dta()
Múltiplos	rio	import()