Crash Course de R

Medidas Resumo e Manipulação de dados

Prof. Carlos Trucíos



Faculdade de Administração e Ciências Contábeis, Universidade Federal de Rio de Janeiro

ctruciosm.github.io — Carlos Trucíos (FACC/UFRJ

Medidas resumo



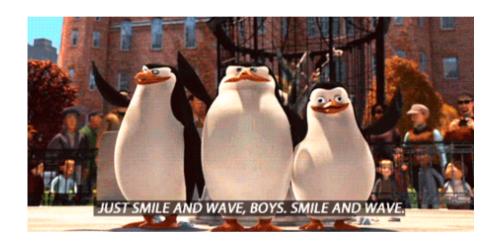
Medidas resumo

```
max(variavel)
                                  # maximo
min(variavel)
                                  # mínimo
mean(variavel)
                                  # média
median(variavel)
                                  # mediana
quantile(variavel, prob = k/100) # k-ésimo percentil
                                  # amplitude inter-quartil
IQR(variavel)
var(variavel)
                                  # variância amostral
sd(variavel)
                                  # desvio padrão amostral
cov(variavel_1, variavel_2)
                                  # covariância
cor(variavel_1, variavel_2)
                                  # correlação de Pearson
summary(dataset_ou_variavel)
                                  # Algumas estatística resumo
table(variavel_categorica)
                                  # Frequências absolutas
prop.table(table(var_categorica)) # Frequências relativas
boxplot(variavel)
                              # Boxplot
hist(variavel)
                              # Histograma
barplot(table(variavel))
                              # Gráfico de barras
plot(variavel_1, variavel_2) # Gráfico de dispersão
```



Medidas resumo

Hands-on:



Utilizaremos o *dataset* penguins do pacote palmerpenguins para fazer uma análise explotarória de dados.

ctruciosm.github.io — Carlos Trucíos (FACC/UFR)

Manipulação de dados





O pacote dplyr

Comando básicos:

- %>%: Pipe
- mutate(): cria novas variáveis.
- select(): seleciona um conjunto de variaveis.
- filter(): filtra casos.
- arrange(): ordena os dados.
- glimpse(): parecido com head()



Importando os dados

```
uri <-"https://raw.githubusercontent.com/ctruciosm/ISLR/master/dataset/Advertising.csv"
dados_advertising <- read.csv(uri)</pre>
```

Carregando o pacote + glimpse()



Suponha que estamos interessados apenas em TV e Sales

```
dados_advertising %>%
    select(TV,Sales)
```

```
TV Sales
##
## 1
       230.1 22.1
## 2
       44.5 10.4
## 3
       17.2
             9.3
       151.5 18.5
## 4
## 5
       180.8 12.9
## 6
        8.7
             7.2
       57.5 11.8
## 7
## 8
       120.2 13.2
## 9
         8.6
              4.8
## 10
      199.8
            10.6
## 11
       66.1
              8.6
## 12
       214.7 17.4
## 13
       23.8
              9.2
## 14
       97.5
              9.7
## 15
      204.1 19.0
## 16
      195.4
             22.4
       67.8 12.5
## 17
## 18
      281.4 24.4
## 19
       69.2 11.3
      147.3 14.6
## 20
## 21
      218.4 18.0
## 22
      237.4 12.5
## 23
       13.2
              5.6
      228.3 15.5
## 24
       62.3
## 25
              9.7
## 26
      262.9 12.0
      142.9 15.0
## 27
      240.1 15.9
## 28
## 29
      248.8 18.9
```



Suponha que estamos interessados em TV, Sales e Sales\$^2\$

```
dados_advertising %>%
  select(TV,Sales) %>%
  mutate(Sales2 = Sales^2)
```

```
TV Sales Sales2
##
      230.1 22.1 488.41
## 1
## 2
       44.5 10.4 108.16
              9.3 86.49
## 3
       17.2
## 4
      151.5 18.5 342.25
      180.8 12.9 166.41
## 5
## 6
        8.7
             7.2 51.84
## 7
       57.5 11.8 139.24
      120.2 13.2 174.24
## 8
## 9
        8.6
              4.8 23.04
## 10
      199.8 10.6 112.36
## 11
       66.1
              8.6 73.96
## 12
      214.7 17.4 302.76
## 13
       23.8
              9.2 84.64
## 14
       97.5
              9.7 94.09
## 15
      204.1 19.0 361.00
## 16
      195.4
             22.4 501.76
       67.8 12.5 156.25
## 17
      281.4 24.4 595.36
## 18
       69.2 11.3 127.69
## 19
## 20
      147.3 14.6 213.16
      218.4 18.0 324.00
## 21
      237.4 12.5 156.25
## 22
       13.2
## 23
              5.6 31.36
      228.3 15.5 240.25
## 24
             9.7 94.09
## 25
       62.3
## 26
      262.9 12.0 144.00
      142.9 15.0 225.00
## 27
## 28
      240.1 15.9 252.81
```



E se quisermos TV, Sales e Sales2 para os valores nos quais as vendas foram (>15)?

```
dados_advertising %>%
  select(TV,Sales) %>%
  mutate(Sales2 = Sales^2) %>%
  filter(Sales>15)
```

```
##
        TV Sales Sales2
## 1 230.1 22.1 488.41
    151.5 18.5 342.25
     214.7 17.4 302.76
     204.1 19.0 361.00
    195.4 22.4 501.76
     281.4
           24.4 595.36
    218.4 18.0 324.00
    228.3 15.5 240.25
## 8
## 9 240.1 15.9 252.81
## 10 248.8 18.9 357.21
## 11 292.9 21.4 457.96
## 12 265.6 17.4 302.76
## 13 266.9 25.4 645.16
## 14 228.0 21.5 462.25
## 15 202.5 16.6 275.56
## 16 177.0 17.1 292.41
## 17 293.6 20.7 428.49
## 18 239.9 23.2 538.24
## 19 216.4 22.6 510.76
## 20 182.6 21.2 449.44
## 21 262.7 20.2 408.04
## 22 198.9 23.7 561.69
## 23 210.8 23.8 566.44
## 24 210.7 18.4 338.56
## 25 261.3 24.2 585.64
## 26 239.3 15.7 246.49
## 27 131.1 18.0 324.00
```



E se quisermos os valores ordenados (de menor a maior) por Sales?

```
dados_advertising %>%
  select(TV,Sales) %>%
  mutate(Sales2 = Sales^2) %>%
  filter(Sales>15) %>%
  arrange(Sales)
```

```
TV Sales Sales2
##
## 1 193.2 15.2 231.04
## 2 123.1 15.2 231.04
      93.9 15.3 234.09
    228.3 15.5 240.25
     184.9 15.5 240.25
    141.3 15.5 240.25
    187.8 15.6 243.36
## 8 239.3 15.7 246.49
     240.1 15.9 252.81
## 9
## 10 209.6 15.9 252.81
## 11 125.7 15.9 252.81
## 12 286.0 15.9 252.81
## 13 110.7 16.0 256.00
## 14 280.7 16.1 259.21
## 15 202.5 16.6 275.56
## 16 197.6 16.6 275.56
## 17 109.8 16.7 278.89
## 18 163.3 16.9 285.61
## 19 213.4 17.0 289.00
## 20 177.0 17.1 292.41
## 21 215.4 17.1 292.41
## 22 135.2 17.2 295.84
## 23 191.1 17.3 299.29
## 24 149.7 17.3 299.29
## 25 214.7 17.4 302.76
## 26 265.6 17.4 302.76
```



E se quisermos ordenados de maior a a menor?

```
dados_advertising %>%
  select(TV,Sales) %>%
  mutate(Sales2 = Sales^2) %>%
  filter(Sales>15) %>%
  arrange(desc(Sales))
```

```
TV Sales Sales2
##
## 1 276.9 27.0 729.00
## 2
     287.6
           26.2 686.44
     283.6 25.5 650.25
    266.9 25.4 645.16
## 4
## 5
     289.7 25.4 645.16
    243.2 25.4 645.16
## 6
## 7 220.3 24.7 610.09
     281.4 24.4 595.36
## 8
     261.3 24.2 585.64
## 9
## 10 210.8 23.8 566.44
## 11 296.4 23.8 566.44
## 12 198.9 23.7 561.69
## 13 239.9 23.2 538.24
## 14 216.4 22.6 510.76
## 15 205.0 22.6 510.76
## 16 195.4 22.4 501.76
## 17 216.8 22.3 497.29
## 18 250.9
           22.2 492.84
## 19 230.1 22.1 488.41
## 20 241.7 21.8 475.24
## 21 213.5 21.7 470.89
## 22 228.0 21.5 462.25
## 23 292.9 21.4 457.96
## 24 182.6 21.2 449.44
## 25 273.7 20.8 432.64
## 26 293.6 20.7 428.49
```



- O %>% nos ajuda a fazer nosso código mais fácil de entender.
- Se formos quebrar o código em váriaslinhas, o %>% deve ir **sempre** no final da linha.

Podemos também calcular algumas estatísticas

```
dados_advertising %>%
  select(TV,Sales) %>%
  mutate(Sales2 = Sales^2) %>%
  filter(Sales>15) %>%
  arrange(desc(Sales)) %>%
  summarise(media_TV = mean(TV), media_Sales = mean(Sales))

## media_TV media_Sales
## 1 213.9013   19.61067
```



Hands-on

Utilize o dataset dados_advertising, filtre os dados para considerarmos unicamente os casos em que Sales <= median(Sales), selecione apenas as variáveis TV, Radio e Newspaper, calcule a media e desvio padrão desses dados.

Gabarito

Cuidado com a ordem dos comandos!



Gabarito

Por quê?



Mais comandos:

- summarise()
- top_n()
- group_by()contains()
- rename()



Podemos calcular estatísticas por grupos.

```
dados_advertising %>%
  select(TV, Sales, Radio, Newspaper) %>%
  group_by(Sales > median(Sales)) %>%
  summarise(media_TV = mean(TV),
            media_Radio = mean(Radio),
            mean_Newspaper = mean(Newspaper))
## # A tibble: 2 × 4
    `Sales > median(Sales)` media_TV media_Radio mean_Newspaper
    <lgl>
                                            <dbl>
                                                           <dbl>
                                <dbl>
## 1 FALSE
                                 94.1
                                             15.5
                                                            27.9
## 2 TRUE
                                             31.3
                                                            33.3
                                202.
```



Se quisermos as 5 lojas com mais vendas?

```
dados_advertising %>%
  select(TV, Sales, Radio, Newspaper) %>%
  top_n(5,Sales)
```

```
## TV Sales Radio Newspaper
## 1 266.9 25.4 43.8 5.0
## 2 289.7 25.4 42.3 51.2
## 3 243.2 25.4 49.0 44.3
## 4 276.9 27.0 48.9 41.8
## 5 287.6 26.2 43.0 71.8
## 6 283.6 25.5 42.0 66.2
```



Se quisermos as 5 lojas com menos vendas?

```
dados_advertising %>%
  select(TV, Sales, Radio, Newspaper) %>%
  top_n(-5,Sales)
```

```
TV Sales Radio Newspaper
## 1 8.6
           4.8
                2.1
                          1.0
## 2 5.4
          5.3 29.9
                          9.4
## 3 13.1
          5.3
               0.4
                         25.6
## 4 0.7
          1.6 39.6
                          8.7
          3.2 11.6
## 5 4.1
                          5.7
```



75.0

7.2

6 48.9

Podemos selecionar variáveis por alguma caracteristica em especial:

```
dados_advertising %>%
    select(contains("a")) %>%
    head()

## Radio Newspaper Sales
## 1 37.8 69.2 22.1
## 2 39.3 45.1 10.4
## 3 45.9 69.3 9.3
## 4 41.3 58.5 18.5
## 5 10.8 58.4 12.9
```



Podemos renomear as variáveis

```
X tv_gastos radio_gastos newspaper_gastos Sales
          230.1
                        37.8
                                        69.2 22.1
## 1 1
## 2 2
           44.5
                        39.3
                                        45.1 10.4
## 3 3
           17.2
                        45.9
                                        69.3 9.3
## 4 4
          151.5
                        41.3
                                        58.5 18.5
## 5 5
          180.8
                                        58.4 12.9
                        10.8
## 6 6
            8.7
                        48.9
                                        75.0 7.2
```



Hands-on

Utilize o dataset dados_advertising e:

- 1. Apage a coluna X (Dica: select(-X))
- 2. Calcule a média (TV, Radio e Newspaper) e o número de elementos por grupo (Sales > mean(Sales)). Dica: para calcular o número de elementos use n().

Gabarito