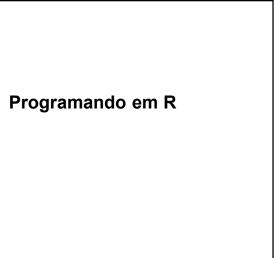
# Estatística Computacional I Lupércio França Bessegato Dep. de Estatística/UFJF



# 1. Programando em R 2. Preparação, limpeza e manipulação de dados 3. Gráficos em R 4. Tópicos especiais 5. Referências

Família Apply



# Família apply



• Podem ser uma alternativa para loops

Função	Descrição		
apply	Aplica funções nas margens de arrays		
by	Aplica uma função em dataframe estratificado por fatores		
eapply	Aplica uma função em valores de um environment		
lapply	Aplica uma função em uma lista, vetor ou data frame		
mapply	Aplica uma função a múltiplos argumentos de lista ou vetor		
rapply	Aplica uma função recursivamente a uma lista		
sapply	Aplica uma função em uma lista, vetor ou data frame		
sweep			
tapply	Aplica uma função sobre uma matriz irregular		
	65		
	Estatística Computacional I - 2020		





- Exemplo Envasamento de leite
  - √ Volume da embalagem ~ Normal com média 1.000 ml e desvio padrão 4 ml
  - √ Amostras aleatórias de tamanho 5, coletadas a cada 30 min, num período de 8 horas

```
> # exemplo -leite
> set.seed(666)
> n <- 5; m <- 16
> amostras <- matrix(rnorm(n*m, 1000, 4), ncol = n, byrow = T)
> amostras <- round(amostras, 2)
> dim(amostras)
[1] 16 5
> rownames(amostras) <- paste(seq(0.5, 8, by = 0.5), "h")
> colnames(amostras) <- paste("Emb.#", 1:5)
> head(amostras)
Emb.# 1 Emb.# 2 Emb.# 3 Emb.# 4 Emb.# 5
0.5 h 1003.01 1008.06 998.58 1008.11 991.13
1 h 1003.03 994.78 996.79 992.83 999.83
1.5 h 1008.60 992.92 1003.46 993.12 1000.54
```

Estatística Computacional I - 2020



# Comando apply



- · Aplica uma função a cada linha ou coluna de uma matriz
- Sintaxe:

```
apply(X, MARGIN, FUN, ...)
\sqrt{X}: um array
\sqrt{\text{MARGIN}}: (1 = linha, 2 = coluna, c(1,2) =
 linha e coluna)
\sqrt{\text{FUN}}: função a ser aplicada
√ . . . : argumentos opcionais da função
```

Estatística Computacional I - 2020

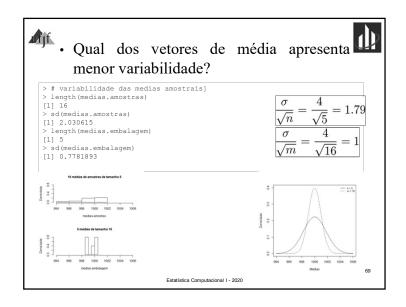


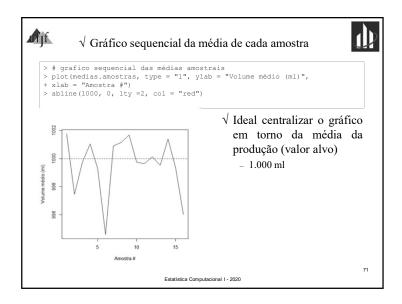
• Exemplo – Caixas de leite

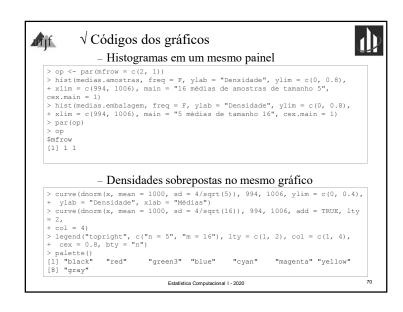


√ Função apply (maneira simples!):

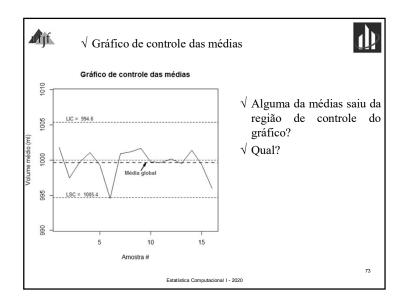
```
> # média por linha - apply (maneira simples)
> medias.amostras <- apply(amostras, 1, mean)
[1] 1001.778 997.452 999.728 1001.064 999.290 994.576 1000.918 1001.158
[9] 1001.690 999.750 999.654 1000.140 999.518 1001.412 999.364 996.022
> is.vector(medias.amostras)
[1] TRUE
> all.equal(medias.amostras, medias.row)
[1] TRUE
> # média por coluna - apply
> medias.embalagem <- apply(amostras, 2, mean)
> medias.embalagem
[1] 999.9094 998.9438 998.6206 1000.4650 1000.0344
```

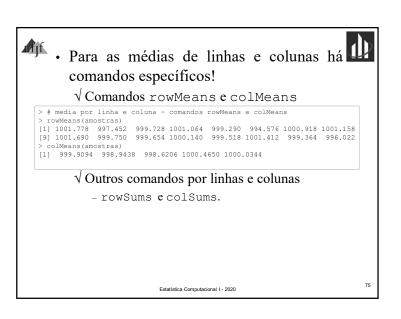


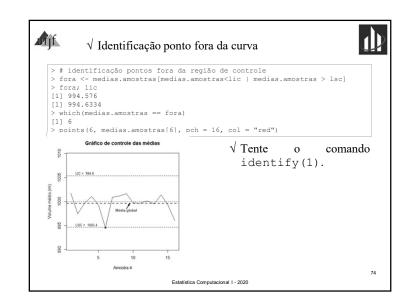


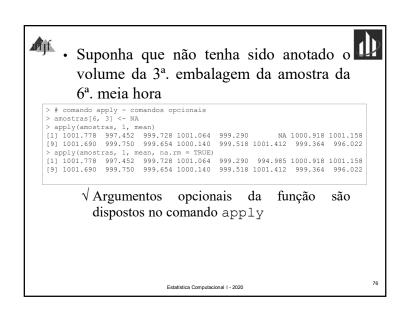


```
Aff.
           √ Gráfico de controle das médias
   > # grafico de controle das médias - a cada 1/2 h
   > lm <- media
   > lsc <- media + 3*dp/sqrt(n)
   > lic <- media - 3*dp/sqrt(n)
   > LIC <- paste("LIC = ", round(lic, 1))
   > LSC <- paste("LSC = ", round(lsc, 1))
   > plot(medias.amostras, type = "1", ylim = c(lic - 4, lsc + 4),
   + ylab = "Volume médio (ml)", xlab = "Amostra #")
   > abline(h = c(lic, 1000, lsc), lty = 2, col = c("blue", "red", "blue"))
   > text(x = c(1, 1), y = c(lic, lsc) + 0.5, c(LSC, LIC), cex = 0.8, pos = 4,
   offset = 1,
   + col = "blue")
   > title(main = "Gráfico de controle das médias")
   > abline(h = mean(medias.amostras), lty = 2, lwd = 2, col = "red")
   > text(9, 1000 - 1, "Média global", pos = 1, cex = 0.8, font = 2, col =
   > arrows (9, 1000 - 1.5, 9.5, mean (medias.amostras), lwd = 2, length = 0.1,
   + angle = 15)
                                                                               72
                                Estatística Computacional I - 2020
```









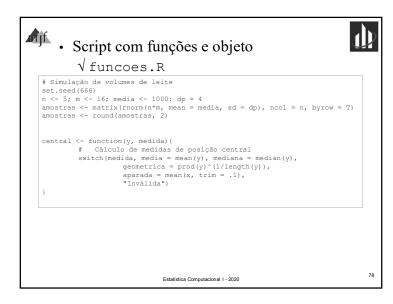


· Deseja-se calcular o coeficiente de variação por linha e por coluna √ Não há comando no R para esse cálculo

> # cálculo de função construída pelo usuário > apply(amostras, 1, function(x) sd(x, na.rm = T)/mean(x, na.rm = T)) [1] 0.0071406716 0.0040634290 0.0067715184 0.0024154519 0.0037076493 [6] 0.0007792574 0.0035788397 0.0050205090 0.0038794314 0.0027465927 [11] 0.0031722062 0.0037257097 0.0055418110 0.0031216872 0.0028533020 [16] 0.0045046186 > apply(amostras, 2, function(x) sd(x, na.rm = T)/mean(x, na.rm = T)) [1] 0.003696620 0.004012343 0.004656646 0.005019586 0.004160150

Estatística Computacional I - 2020







• Aplicação de função do usuário

√Função central, para cálculo de médias geométricas por linha

> # aplicação de função construída anteriormente > source("funcoes.R") > apply(amostras, 1, FUN = central, medida = "geometrica") [1] 1001.7575 997.4454 999.7097 1001.0617 999.2845 994.5755 1000.9129 [8] 1001.1479 1001.6840 999.7470 999.6500 1000.1345 999.5057 1001.4081 [15] 999.3607 996.0139

Estatística Computacional I - 2020



• apply:



80

√ Aplica função a todos os elementos de objeto

Função	Objetos(a)	Elementos(b)	Tipo de resultado
apply	Matriz	Linhas ou colunas	Vetor, matriz, array ou lista
	Array	Linhas, colunas ou qualquer dimensão	Vetor, matriz, array ou lista
	Data frame		Vetor, matriz, array ou lista

 $\sqrt{}$  (a) Objeto com os quais a função trabalha

 $\sqrt{(b)}$  Elementos do objeto vistos pela função



# Comando lapply



- · Aplica uma função a cada linha ou coluna de uma matriz
  - √ Saída é uma lista
- Sintaxe:

```
lapply(X, FUN, ...)
```

- $\sqrt{X}$ : vetor, lista ou data frame
- √ FUN: função a ser aplicada
- √ . . . : argumentos opcionais da função

Estatística Computacional I - 2020





83

- Diferença com o comando apply:
  - √Pode ser usada com objetos do tipo data frame, lista ou vetor
  - √ A saída do comando é uma lista
    - Possui o mesmo número de elementos do objeto

Estatística Computacional I - 2020



# • lapply:

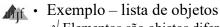


√ Aplica função a todos os elementos de objeto

Função	Objetos(a)	Elementos(b)	Tipo de resultado
lapply	Vetor	Elementos	Lista
	Data frame	Variáveis	Lista
	Lista	Elementos	Lista

- $\sqrt{}$  (a) Objeto com os quais a função trabalha
- $\sqrt{(b)}$  Elementos do objeto vistos pela função

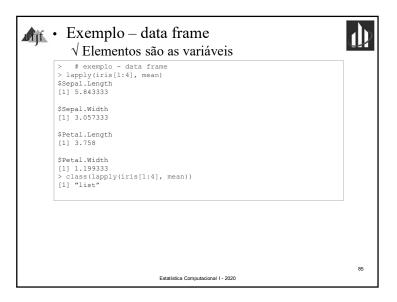
Estatística Computacional I - 2020

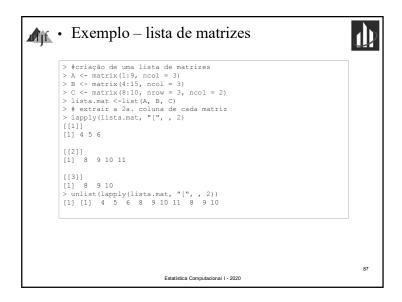


√ Elementos são objetos diferentes

```
> #cria lista com dois elementos
> lista <- list(a = 1:10, b = 11:20, c = matrix(1:6, ncol = 2))
> # média dos valores em cada elemento da lista
> lapply(lista, mean)
$a
[1] 5.5
[1] 15.5
[1] 3.5
> class(lapply(lista, mean))
[1] "list"
> # soma dos valores em cada elemento
> lapply(lista, sum)
[1] 55
[1] 155
$c
[1] 21
                         Estatística Computacional I - 2020
```

Prof. Lupércio F. Bessegato - UFJF







# • Exemplo – vetor

√ Função é aplicada em cada elemento do vetor

```
> # exemplo - vetor
> lapply(1:3, function(x) x^2)
[[1]]
[[1]]

[[2]]
[1] 4

[[3]]
[1] 9
> # transforma lista em vetor
> unlist(lapply(1:3, function(x) x^2))
[1] 1 4 9
```

• Importante:

√ Se você quer uma lista como saída, use lapply; se quer um vetor, use sapply

Estatística Computacional I - 2020

**A**jf.

## Comando sapply



88

 Aplica uma função a cada linha ou coluna de uma matriz

√ Saída é um vetor

• Sintaxe:

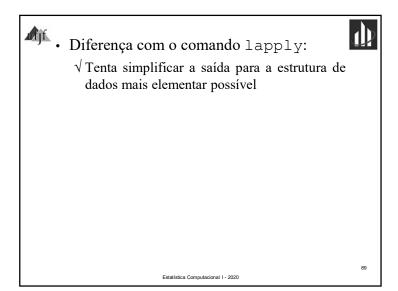
```
sapply(X, FUN, ..., simplify = T)

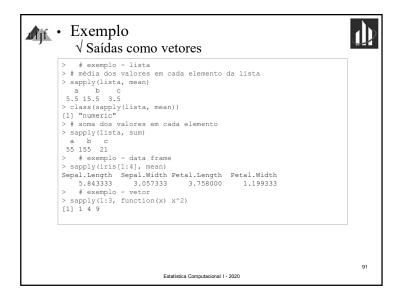
√X: vetor, lista ou data frame

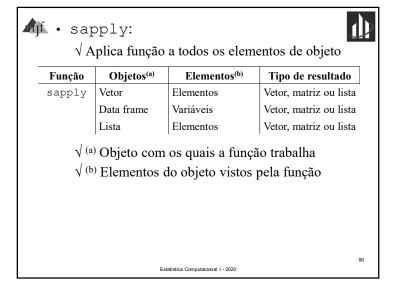
√FUN: função a ser aplicada

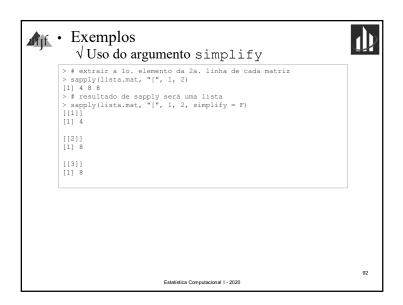
√...: argumentos opcionais da função

√ simplify: T = saída é vetor, F = saída é lista
```











# • Exemplo:



√Rendimento total de cevada, em bushels por acre, para 10 variedades em 6 locais, em 2 anos

```
> # barley data frame
> library (lattice)
> dim (barley)
[1] 120 4
> head(barley)
   yield variety year
1 27.00000 Manchuria 1931 University Farm
2 48.86667 Manchuria 1931
3 27.43334 Manchuria 1931
4 39.93333 Manchuria 1931
                              Crookston
5 32.96667 Manchuria 1931 Grand Rapids
6 28 96667 Manchuria 1931
                                 Duluth
> unique(barley$site)
[1] University Farm Waseca
                                                  Crookston
[5] Grand Rapids Duluth
> unique(barley$yield)
> barley$yield[duplicated(barley$yield)]
[1] 29.66667 48.56666 32.96667 35.13333 27.43334 20.63333
```



# √ Quantidade de valores únicos por variável



```
> # valores únicos por variável
> # saída: lista
> lapply (barley, function(x) length(unique(x)))
$yield
[1] 114
Svariety
[1] 10
$year
[1] 2
Ssite
> is.list(lapply (barley, function(x) length(unique(x))))
> # saída: vetor
> sapply (barley, function(x) length(unique(x)))
 yield variety year site
   114 10
> # funciona em data frames mas não em listas
> apply (barley, 2, function(x) length(unique(x)))
 yield variety year site
   114 10
                   2
> is.vector(apply (barley, 2, function(x) length(unique(x))))
```

Estatística Computacional I - 2020



# Comando tapply

Estatística Computacional I - 2020



95

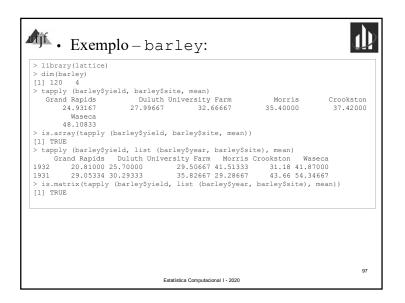
- Aplica uma função a estratos de um vetor √ Saída é um vetor
- Sintaxe:

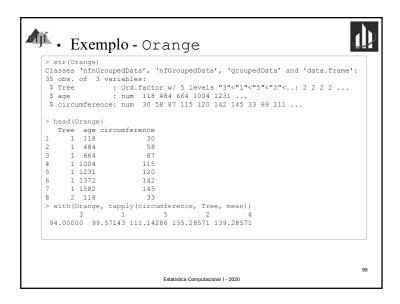
```
tapply(X,
                    INDEX,
                                 FUN.
  simplify = T)
\sqrt{X}: vetor, lista ou data frame
√ INDEX: lista de um ou mais fatores, com
  mesma dimensão que X.
√ FUN: função a ser aplicada
√ . . . : argumentos opcionais da função
\sqrt{\text{simplify: F} = \text{saida \'e lista}}
```

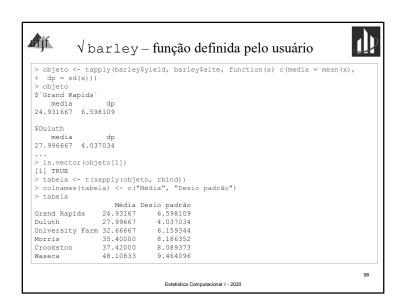
Estatística Computacional I - 2020

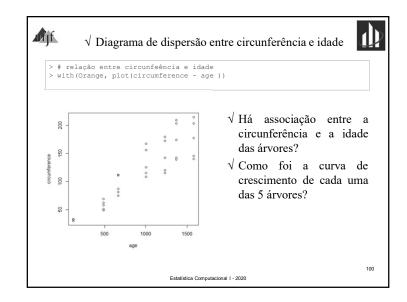
# • Exemplo - mtcars

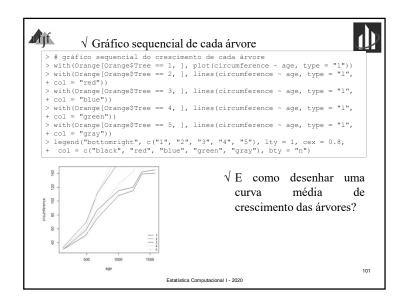
```
num [1:32] 6 6 4 6 8 6 8 4 4 6 ...
> is.factor(mtcars$cyl)
[1] FALSE
> levels(as.factor(mtcars$cyl))
[1] "4" "6" "8"
> tapply(mtcars$mpg,mtcars$cyl,mean)
26.66364 19.74286 15.10000
> # media de mpg por nivel das variáveis cyl e gear
> with (mtcars, tapply (mpg, list (cyl, gear), mean, na.rm = T))
4 21.50 26.925 28.2
6 19.75 19.750 19.7
8 15.05 NA 15.4
```

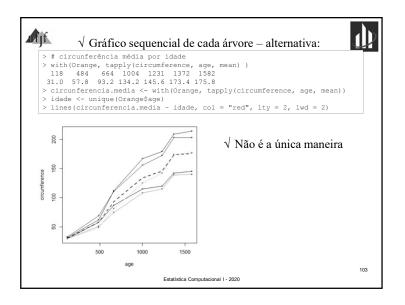


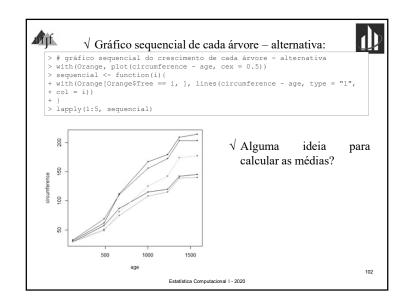


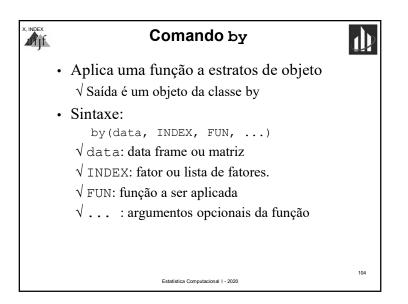


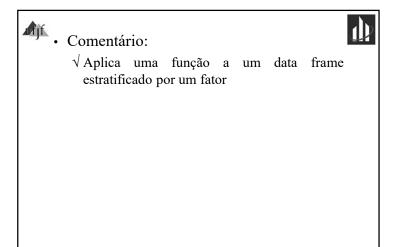




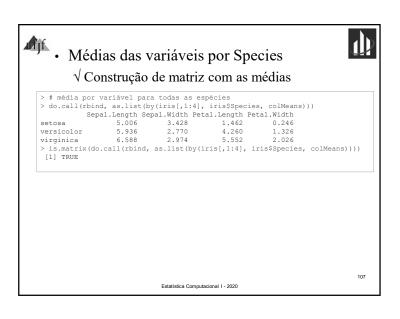


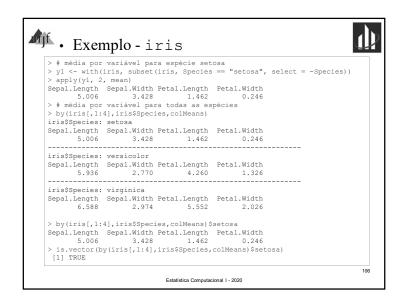






Estatística Computacional I - 2020





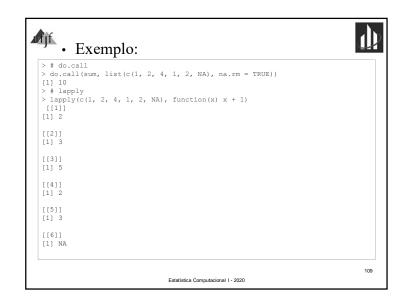


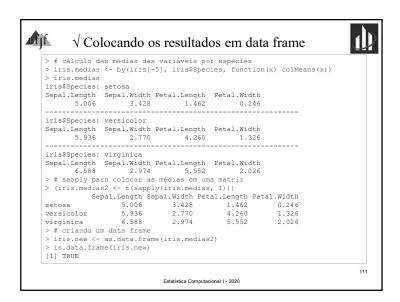
### Comando do . call

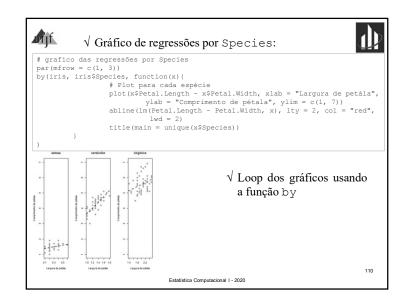


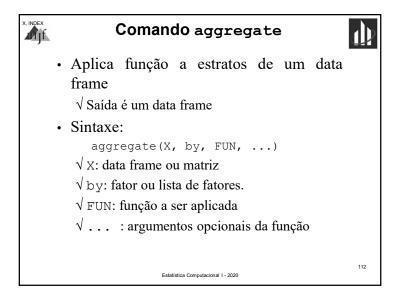
108

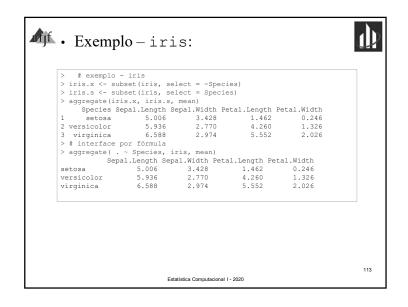
- Executa uma função cujos argumentos estão em uma lista
  - √ Diferente de lapply (ou sapply)
- lapply:
  - √ Aplica uma função a todos os elementos de uma lista

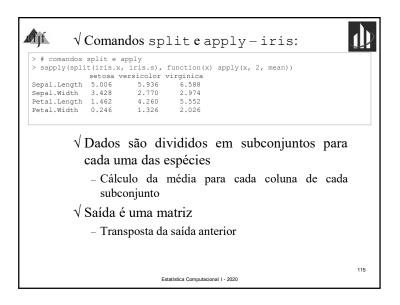






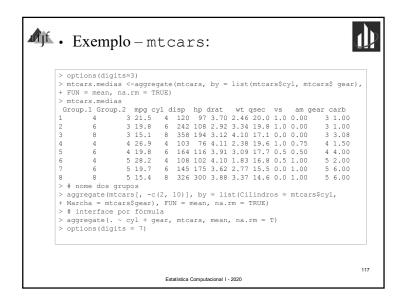


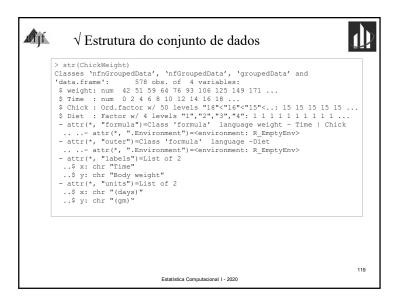


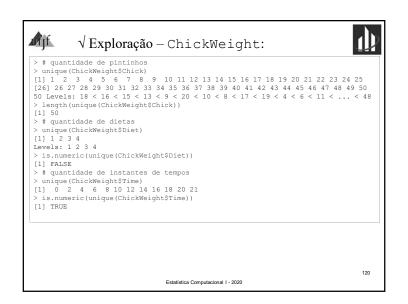


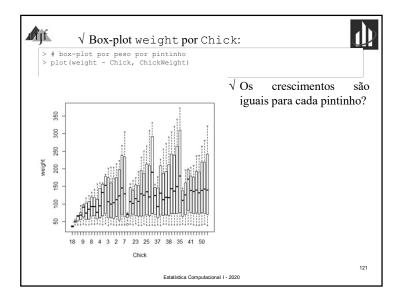
```
√ Média e desvio padrão – iris:
> # função definida pelo usuário
> # calculo de mais de uma função por estrato
> ag <- aggregate(. ~ Species, iris, function(x) c(mean = mean(x), sd = sd(x)))
> is.data.frame(ag)
[1] TRUE
> dim(ag)
[1] 3 5
     Species Sepal.Length.mean Sepal.Length.sd Sepal.Width.mean Sepal.Width.sd
                    5.0060000 0.3524897 3.4280000
5.9360000 0.5161711 2.7700000
2 versicolor
                                                                  0.3137983
                   6.5880000
                                 0.6358796
                                                   2.9740000
                                                                  0.3224966
3 virginica
  Petal.Length.mean Petal.Length.sd Petal.Width.mean Petal.Width.sd
        1.4620000 0.1736640 0.2460000 0.1053856
          4.2600000
                         0.4699110
                                        1.3260000
                                                       0.1977527
         5.5520000
                        0.5518947
                                        2.0260000
> ag[[1]]
[1] setosa
             versicolor virginica
Levels: setosa versicolor virginica
> ag[-1][[4]]
 mean
[1,] 0.246 0.1053856
[2,] 1.326 0.1977527
[3,] 2.026 0.2746501
                                                                          114
                              Estatística Computacional I - 2020
```

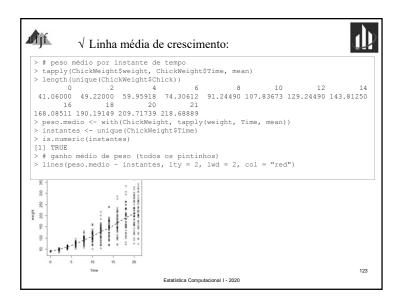
```
• Exemplo – datas:
> datas <- data.frame(data = as.Date("01-01-2018", format = "%d-%m-%Y") + 0:729)
 > head(datas)
 data
 1 2018-01-01
 2 2018-01-02
3 2018-01-03
 > mean (datas$data)
 [1] "2018-12-31"
 > as.Date("28/08/1991", format = "%d/%m/%Y") - as.Date("17/01/1991",
 + format = "%d/%m/%Y")
Time difference of 223 days
 > as.Date("31/07/2017", format = "%d/%m/%Y") + 223
 [1] "2018-03-11"
 > # último dia de cada mês
 > ultimo <- aggregate(x = datas["data"],
                       by = list(mes = substr(datas$data, 1, 7)),
                        FUN = mean)
> head(ultimo)
                                         > substr("abcdef", 2, 4)
 1 2018-01 2018-01-16
 2 2018-02 2018-02-14
                                         [1] "bcd"
 3 2018-03 2018-03-16
                                                                             116
                               Estatística Computacional I - 2020
```

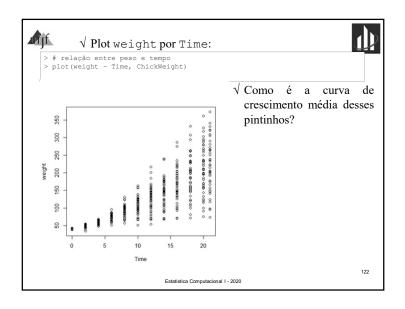


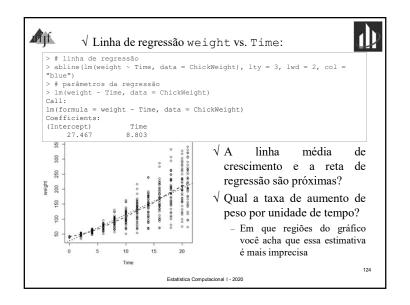




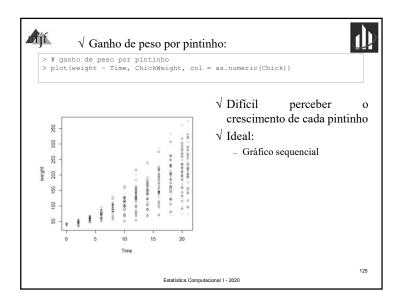


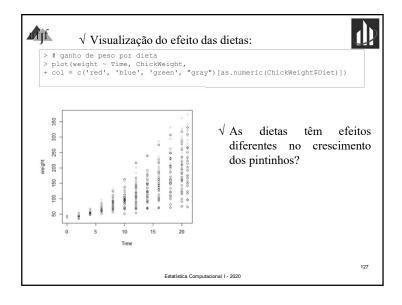


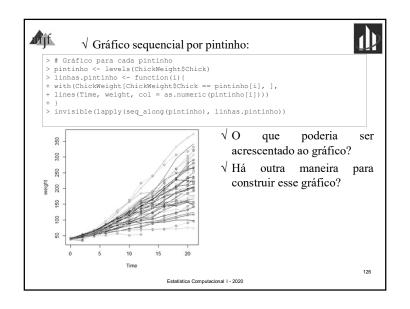




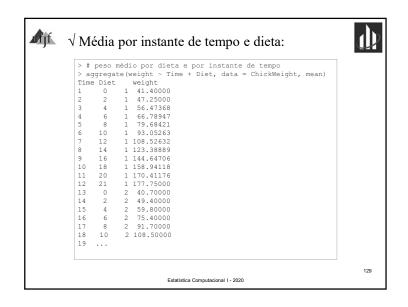
Prof. Lupércio F. Bessegato - UFJF

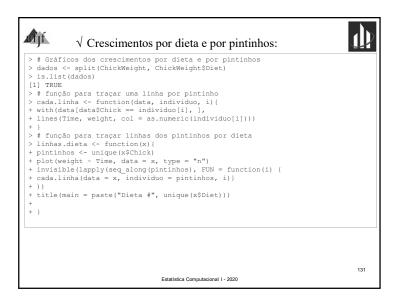


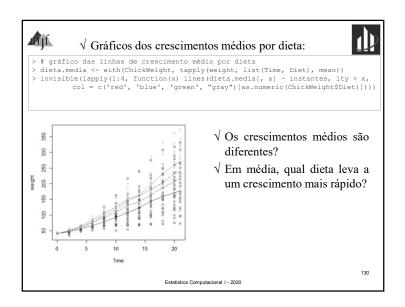


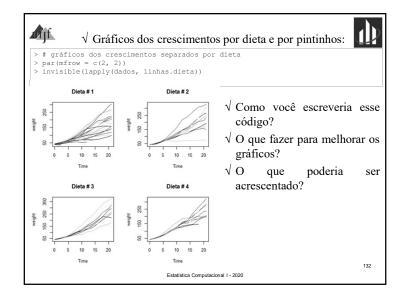


```
\sqrt{\text{Estatísticas por dieta tempo:}}
> # média de peso por dieta
> with (ChickWeight, aggregate (weight, list(dieta = Diet), mean))
  dieta
     1 102.6455
2 122.6167
      3 142.9500
     4 135.2627
> # desvio padrão por dieta
> with(ChickWeight, aggregate(weight, list(dieta = Diet), sd))
  dieta
      1 56.65655
      2 71.60749
      3 86.54176
     4 68.82871
> #peso médio por instante de tempo
 > aggregate(weight ~ Time, data = ChickWeight, mean)
   Time weight
      0 41.06000
      2 49 22000
      4 59.95918
      6 74.30612
      8 91.24490
     10 107.83673
     12 129.24490
     14 143.81250
                                Estatística Computacional I - 2020
```











# Comando rapply



133

- Aplica recursivamente função a todos os elementos de uma lista
  - √ Saída é um vetor ou uma lista
- Sintaxe:

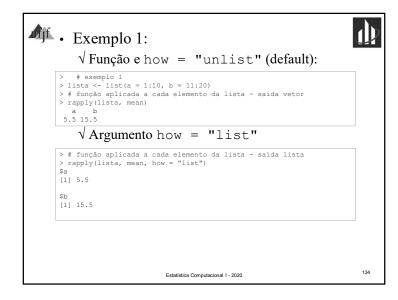
```
rapply(object, f, ...)

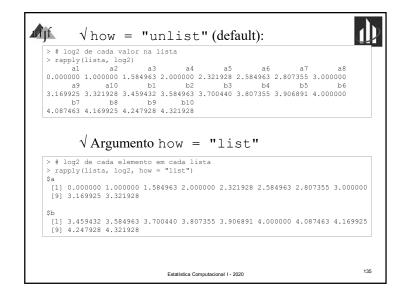
√object: lista

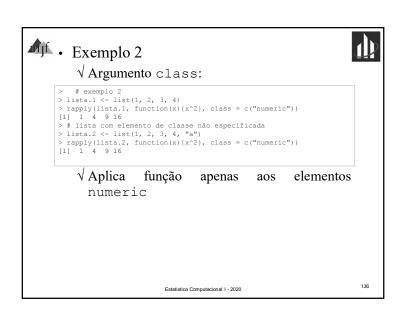
√f: função com um único argumento

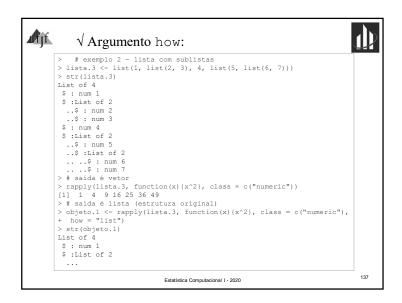
√class: classes a que função deve ser aplicada

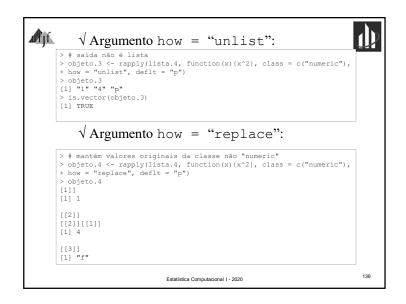
√...: argumentos opcionais da função
```

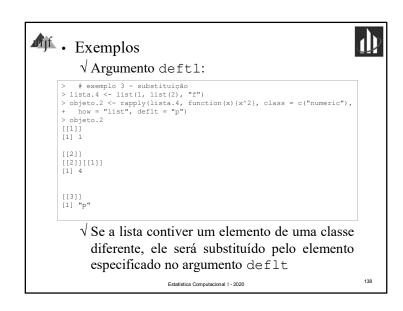


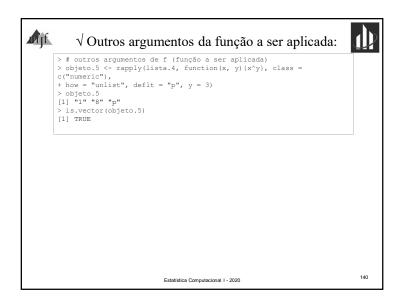














# Comando vapply



- Similar a sapply,
  - √Requer especificação do tipo de dados desejados na saída
  - √ Saída é um vetor ou matriz
- Sintaxe:

```
vapply(X, FUN, FUN.VALUE, ...)
\sqrt{X}: vetor ou lista
√ FUN: função com um único argumento
√ FUN. VALUE: tipo de dados desejados
√ . . . : argumentos opcionais da função
```

Estatística Computacional I - 2020

141



Min. -1.4534842 -0.8821279 1o. Qu. -0.6873098 -0.2106322 Mediana -0.1527475 0.6392059 3o. Qu. 0.4422310 0.9666165 1.4889964 1.8130571

#### √ Saída é uma matriz

- Nomes das colunas: elementos da lista original
- Nomes das linhas: template da saída

Estatística Computacional I - 2020

143



# • Exemplo 1:



#### √ Vetores e listas:

```
> # exemplo 1
> # saida numerica - vetor
> vetor.v <- 1:10
> vapply(vetor.v, sum, numeric(1))
[1] 1 2 3 4 5 6 7 8 9 10
> # saída numérica - lista
> A <- 1:9
> B <- 1:12
> C <- 1.15
> lista.v <- list(A, B, C)
> vapply(lista.v, sum, numeric(1))
[1] 45 78 120
```

#### √ Saída será um erro se função retornar mais de um valor numérico

- FUN. VALUE = numeric (1) pode ser útil caso se espere apenas um resultado elemento

Estatística Computacional I - 2020

142



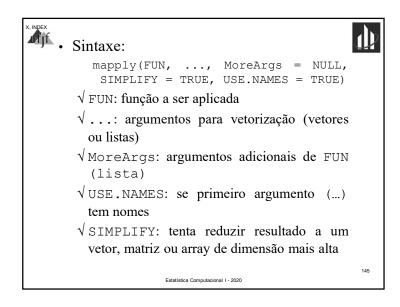
# Comando mapply

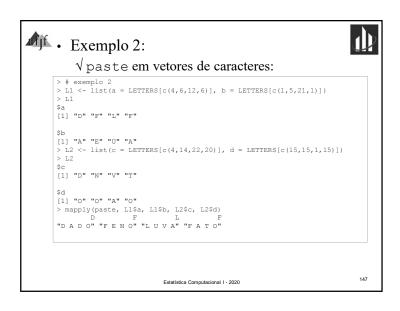


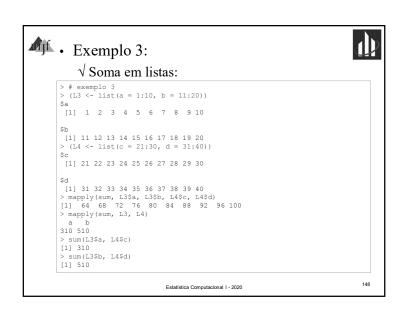
- Versão multivariada de lapply sapply
  - $\sqrt{\text{Caso lapply e sapply:}}$ 
    - Função atua somente sobre os elementos de uma única lista
  - √ Caso mapply
    - Função é aplicada sobre o 1º elemento de cada um dos argumentos, em seguida ao 2º, etc.
    - Argumentos podem ser listas ou vetores
    - Saída é um vetor ou matriz

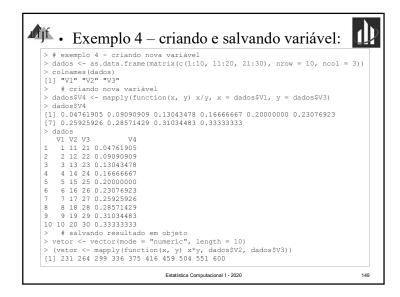
Estatística Computacional I - 2020

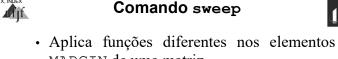
144

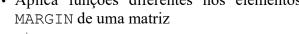










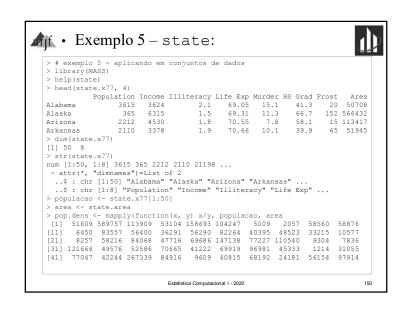


√ Aplicação de função em linha (coluna) de matriz, quando o outro argumento da função tem valor diferente para cada linha (coluna)

√ Saída é um vetor ou matriz

Estatística Computacional I - 2020

151



```
• Sintaxe:

sweep (x, MARGIN, STATS, FUN="-
",...)

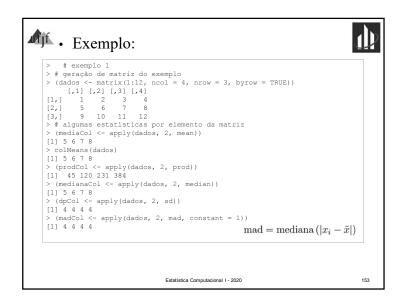
√x: array

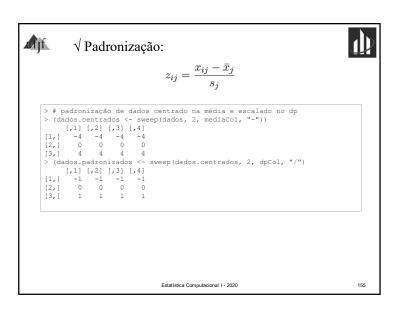
√MARGIN: 1 = linha, 2 = coluna

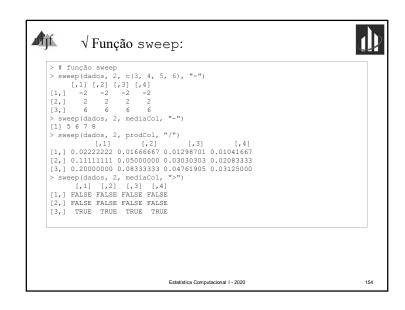
√STATS: estatística do elemento escolhido

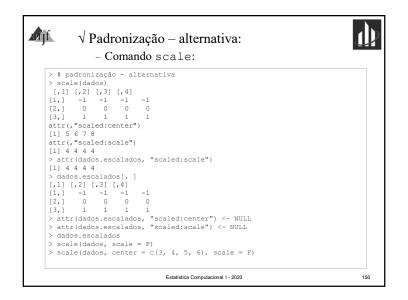
√FUN: função que aplicará STATS

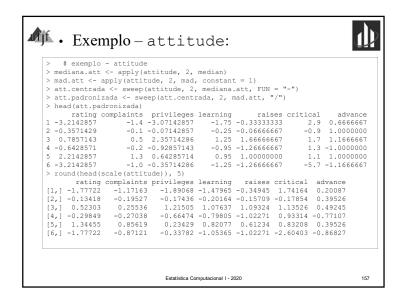
√...: argumentos opcionais da função
```



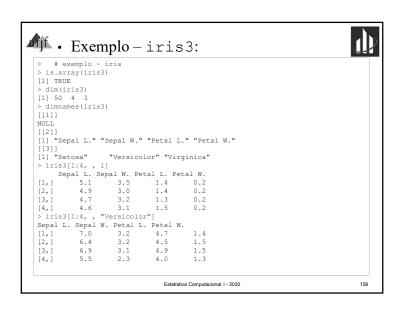




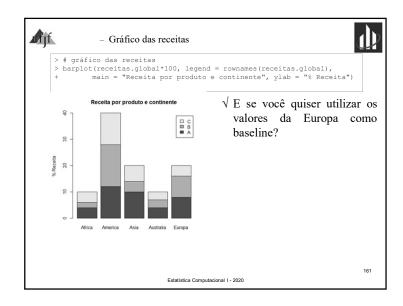




```
Aff.
         √ Família apply – iris3:
 > # comandos apply e sweep
 > iris.medias <- apply(iris3, c(2, 3), mean)
 > iris.medias
          Setosa Versicolor Virginica
 Sepal L. 5.006 5.936 6.588
 Sepal W. 3.428
                    2.770
                             2.974
  Petal L. 1.462 4.260 5.552
 Petal W. 0.246
                   1.326 2.026
 > iris3.centrada <- sweep(iris3, c(2,3), iris.medias, "-")
 > iris3.centrada[1:2, , ]
 , , Setosa
     Sepal L. Sepal W. Petal L. Petal W.
 [1,] 0.094 0.072 -0.062 -0.046
 [2,] -0.106 -0.428 -0.062 -0.046
 , , Versicolor
      Sepal L. Sepal W. Petal L. Petal W.
 [1,] 1.064 0.43 0.44 0.074 [2,] 0.464 0.43 0.24 0.174
 , , Virginica
      Sepal L. Sepal W. Petal L. Petal W.
 [1,] -0.288 0.326 0.448 0.474
 [2,] -0.788 -0.274 -0.452 -0.126
                             Estatística Computacional I - 2020
                                                                        159
```



```
• Receitas globais de produtos:
 > # exemplo - Receita por produto e continente
> # montagem do conjunto de dados
> produto <- c("A", "B", "C", "Total")
 > continente <- c("Africa", "America", "Asia", "Australia", "Europa")
 > valores <- c(0.4, 0.2, 0.4, 0.1, 0.3, 0.4, 0.3, 0.4, 0.5, 0.2,
             0.3, 0.2, 0.4, 0.3, 0.3, 0.1, 0.4, 0.4, 0.2, 0.2)
 > (receitas <- matrix(valores, ncol = 5, dimnames = list(produto, continente)))
      Africa America Asia Australia Europa
       0.4 0.3 0.5
                            0 4 0 4
        0.2 0.4 0.2
                             0.3 0.4
               0.3 0.3
                             0.3
 Total 0.1 0.4 0.2 0.1 0.2
 > # receitas absolutas
 > (receitas.global <- sweep(receitas[1:3, ], 2, receitas[4, ], "*"))
 Africa America Asia Australia Europa
 A 0.04 0.12 0.10
                        0.04 0.08
 B 0.02 0.16 0.04
                        0.03 0.08
 C 0.04 0.12 0.06
                        0.03 0.04
 > sum(receitas.global)
 [1] 1
                            Estatística Computacional I - 2020
                                                                      160
```





# Bibliografia Recomendada



286

- ALBERT, J.; RIZZO, M. R by Example. Springer, 2012.
- CHRISTIAN, N. Basic Programming, Lecture Notes
- DALGAARD, P. *Introductory statistics with R.* Springer, 2008.
- KLEIBER, C.; ZEILEIS, A. Applied econometrics with R. Springer, 2008.
- GARDENER, M. Beginning R: The statistical programming language. John Wiley & Sons, 2012.

Estatística Computacional I - 2020

Referências