

Introdução à Análise de Dados com Linguagem R

1. Introdução

R é uma linguagem de programação gratuita de código aberto usada principalmente para computação estatística e gráficos. R é uma linguagem interpretada semelhante a Python, onde você não precisa compilar primeiro para executar seu programa. Depois de criar seu programa, você pode executá-lo em uma ampla variedade de plataformas UNIX, Windows e macOS.

R é uma linguagem de programação orientada a objetos, criada em 1996 por Ross Ihaka e Robert Gentleman. R é uma linguagem específica de domínio de código aberto, explicitamente projetada para ciência de dados. Muito popular em finanças e academia, R é uma linguagem perfeita para manipulação, processamento e visualização de dados, bem como computação estatística e aprendizado de máquina.

Como qualquer outra linguagem de programação, R também suporta extensão na forma de pacotes, portanto, os desenvolvedores podem criar seus próprios pacotes e reutilizá-los quando necessário.

2 Instalação do R

A instalação padrão da linguagem R é feita a partir do [CRAN](https://cran.r-project.org/), uma rede formada de servidores espalhados pelo mundo que armazena versões atualizadas do código fonte e executável (Windows), assim como a documentação da linguagem R.

3 Instalação do RStudio

RStudio é um ambiente de desenvolvimento integrado para linguagem R (IDE - *integrated development environment*), porém pode rodar *scripts* SQL, C, C++ e Python. A vantagem de se poder trabalhar com IDE é que ela disponibiliza ferramentas de apoio ao desenvolvimento de códigos em linguagem de programação. Para *download* do RStudio acesse o endereço <https://www.rstudio.com/>, e escolha a opção *RStudio Desktop*

3.1 Conhecendo a Interface Gráfica do RStudio

A interface do RStudio é dividida em 4 painéis, e duas barras:

- 1 - Barra de *menu*
- 2 - Barra de ferramentas
- 3 - Painel de *scripts* e arquivos
- 4 - Painel de variáveis de Ambiente/Histórico/Conexões/Tutorial
- 5 - Painel de *Console/Terminal*
- 6 - Painel de Árvore de Arquivos/Gráficos/Pacotes/Ajuda/Visualizador

4. Obter e Configurar o ambiente de trabalho

Para obter o ambiente de trabalho atualmente em uso pelo RStudio utilizamos a função `getwd()`; esta função vem do termo em inglês: *Get Working Directory*, traduzido para o português como "Obter Ambiente de Trabalho". No sistema operacional Windows, por padrão, o RStudio configura o ambiente de trabalho em "C:/Usuários/Nome_do_Usuário/Documentos"

.

Para configurar um ambiente de trabalho diferente do padrão utilizamos a função `setwd()`. Esta função tem nome bem sugestivo na língua inglesa, a saber: *Set Working Directory*, "Configurar Ambiente de Trabalho".

In []:

Exercício Prático - Ambiente de trabalho

Instrução 1/2

- Use a janela console do RStudio para obter o atual diretório de trabalho em uso no seu computador.

In []:

Instrução 2/2

- Configure seu ambiente de trabalho para `'C:/Users/Downloads'`

In []:

4.1 Operadores Aritméticos e de Atribuição em R

Operador	Função
+	Soma
-	Subtração
/	Divisão
*	Multiplificação
%%	Resto da divisão
%/%	Parte inteira divisão
^	Potenciação
**	Potenciação
<-	Atribuição
=	Atribuição

Exercício Prático - Operadores Aritméticos e de Atribuição

Instrução 1/3

- Obter o resto da divisão entre os números inteiros 10 e 3.

In []:

Instrução 2/3

- Obter a parte inteira da divisão de 10 por 3.

In []:

Instrução 3/3

- Obter o quadrado de um número inteiro qualquer.

In []:

5 Operadores de Comparação em R

Operador	Significado
==	igual a
!=	diferente de
>	maior que
<	menor que
>=	maior ou igual a
<=	menor ou igual a

6 Objeto em R

Um objeto é simplesmente qualquer variável que armazene um caractere numérico, alfabético ou uma cadeia de caracteres (*string*). Em R temos objetos especiais para manipulação de grande volume de dados, a exemplo de vetores, listas e dataframes. Para criação de um objeto se utiliza o operador de atribuição `<-`. Nos exemplos a seguir são criados os objetos `x`, `y` e `z` para armazenar um número inteiro, uma letra e uma frase.

In []:

```
# definição de um objeto do tipo inteiro
int <- 42

# definição de um objeto do tipo caractere
letra <- "R"

# # definição de um objeto do tipo string (cadeia de caracteres)
string <- "R é massa"
```

Em ambos os trechos de código acima se observa o uso do caractere ``#``, em R e python este caractere é utilizado para fazer comentário no código, portanto, toda linha que contém este caractere é ignorada na execução do programa.

A seguir é mostrado como fazer comentário com várias linhas.

In []:

```
"
Este é um comentário em R utilizando
várias linhas para documentação de
códigos.
"
```

6.1 Conferindo o conteúdo de um objeto

Para conferir o conteúdo de um objeto em R, fazemos uma chamada diretamente pelo nome do objeto de interesse ou através do uso da função `print()`

```
In [ ]: int
        print(int)
        cat(int)
```

6.2 Descobrimo o tipo de dados armazenados em um objeto R

Para descobrir o tipo de dados armazenado em um objeto, podemos utilizar a função `class`

```
In [ ]: class(int)
        class(letra)
        class(string)
```

7 Estrutura de Dados em R

7.1 Vetores

Vetor em R é um objeto R que armazena um ou mais elementos de valores indexados, ou seja, cada elemento dentro do vetor possui uma posição específica. Para criação de um vetor basta inserir os valores como argumento dentro da função `c()`, "c" é uma abreviação para a palavra inglesa *concatenate*. Vetor é uma estrutura de dados especialmente importante em análise de dados. A seguir temos um exemplo de como criar um vetor de inteiros.

```
In [ ]: inteiros <- c(42, 33, 0, -1, 5)
```

Para acessar o primeiro elemento do vetor `inteiros` usamos o comando `vetor[x]`, onde *vetor* é nome atribuído ao vetor e *x* é o índice do elemento a ser buscado no vetor. Se quisermos mostrar apenas o elemento de índice 1 do vetor `inteiros`, ou seja o primeiro elemento, basta usar `inteiros[1]`

```
In [ ]:
```

7.1.2 Substituição de elementos de um vetor

```
In [ ]: inteiros[1] <- 2
        inteiros
```

7.2 Funções básicas aplicadas a vetores

- `length()` Retorna o tamanho de um vetor, ou seja, o número de elementos armazenados no vetor.

```
In [ ]: length(inteiros)
```

- `which.max()` Retorna o índice relativo ao maior valor presente nos dados.

In []:

- `which.min()` Retorna o índice relativo ao menor valor presente nos dados.

In []:

- `names()` Retorna os nomes atribuídos a cada elemento de um vetor

In []:

```
dias_semana <- c(1:7)
names(dias_semana)
```

In [41]:

```
names(dias_semana) <- c('Domingo', 'Segunda', 'Terça', 'Quarta', 'Quinta', 'Sexta', 'Sáb
```

In [42]:

```
names(dias_semana)
```

'Domingo' · 'Segunda' · 'Terça' · 'Quarta' · 'Quinta' · 'Sexta' · 'Sábado'

In []:

```
dias_semana
```

- `attributes()` Retorna uma lista com os atributos associados a um vetor.

In []:

```
attributes(inteiros)
attributes(dias_semana)
```

- `seq()` Cria uma sequência dentro de um vetor

In []:

```
sequencia <- seq(0, 10, 2)
print(sequencia)
```

- `rep()` Cria uma repetição de um vetor

In []:

```
rep(1:4, 3)
```

- `duplicated()` Remove elementos duplicados.

In []:

- `unique()` Retorna apenas os valores distintos.

In [68]:

```
# Leitura de vetores de um inventário florestal
load('./data/dados_modulo_1.rda')

# Mostrar os objetos atualmente disponíveis no ambiente R
ls()
```

'altura' · 'categoria' · 'dap' · 'dias_semana' · 'int' · 'inteiros' · 'letra' · 'nomes_cientificos' · 'sequencia' · 'string' ·
'x' · 'y' · 'z'

In []:

7.3 Funções Estatísticas Aplicadas a Vetores

- `mean()` Retorna a média aritmética
- `median()` Retorna a mediana
- `min()` Retorna o menor valor
- `max()` Retorna o maior valor
- `sd()` Retorna o desvio padrão
- `summary()` Retorna a estatística descritiva.
- `cor()` Retorna a correlação entre dois vetores.

Calcular a média do vetor *dap* (diâmetro médio das árvores)

```
In [ ]: # média do vetor dap (diâmetro médio das árvores)
mean(dap)
```

Calcular a mediana do vetor *altura* (altura comercial das árvores)

```
In [ ]: ### Calcular a mediana do vetor _altura_ (altura comercial das árvores)
median(altura)
```

Mostrar os valores mínimo e máximo do vetor *dap*

```
In [ ]: # valor mínimo de dap
min(dap)
```

```
In [ ]: # dap Máximo
max(dap)
```

Calcular o desvio padrão do vetor *dap*

```
In [ ]: # desvio padrão para o vetor dap
sd(dap)
```

Mostrar a estatística Descritiva do Vetor *altura*

```
In [ ]: summary(altura)
```

Calcular a correlação linear entre diâmetro e altura das árvores

```
In [ ]: cor(dap, altura)
```

7.4 Função `lapply()` aplicada a vetores

A função `lapply`, parte do pacote base do R, no caso específico de vetores, recebe 2 argumentos como parâmetro: o vetor contendo de dados e uma função a ser aplicada aos elementos do vetor.

```
In [73]: nomes <- c('MASSARANDUBA', 'IPÊ', 'GARAPEIRA', 'JATOBÁ')
          nomes
```

```
'MASSARANDUBA' · 'IPÊ' · 'GARAPEIRA' · 'JATOBÁ'
```

```
In [ ]: lapply(nomes, tolower)
```

7.5 Função `sapply()`

```
In [ ]: sapply(nomes, tolower)
```

7.6 Função `mapply()`

Versão multivariada das funções `lapply` e `sapply`, utilizada para iterar entre elementos de vetores ou listas.

```
In [ ]: # Definição dos Vetores a e b
a <- c(7, 12, 5, 2, 1)
b <- c(4, 2, 3, 5, 1)

# Nomes para os vetores
dias_semana <- c('Segunda', 'Terça', 'Quarta', 'Quinta', 'Sexta')

# Atribuir nome aos vetores
names(a) <- dias_semana
names(b) <- dias_semana

# Uso da função mapply() para retornar a soma
# entre os elementos dos vetores a e b
mapply(max, a, b)
```

7.7 Função `tapply()`

Aplica uma função sobre um vetor com agrupamento em outro vetor categórico. Recebe como parâmetros: um vetor numérico, um vetor categórico e uma função. O código a seguir aplica a função média sobre o vetor `volume` agrupado ao vetor `ut` (unidades de trabalho)

```
In [ ]: # Calcular o volume médio por unidade de trabalho
tapply(dap, categoria, mean)
```

8 Valores Ausentes (NA)

Em R valores ausentes são conhecidos como `NA`, uma sigla em inglês que significa *Not Available*, ou seja, valores não disponíveis. Na literatura técnica e também em outras linguagens de programação esta sigla é definida como `Nan` (*Not available number*)

In []:

In []:

Como Saber se Há Valores Ausentes (NA) nos Dados?

Para conferir a presença de valores NA nos dados utilizamos a função `is.na()`, a qual recebe como parâmetro de entrada apenas o vetor de dados.

In []:

Exercício Prático - Vetores e Operadores de Comparação

Para este exercício, considere o vetor `temperatura`. Esse vetor possui dados de temperatura média mensal da Estação Meteorológica Manual INMET 82861, localizada no município de Conceição do Araguaia.

```
In [ ]: temperatura <- c(26.38452, 26.90357, 27.04064, 27.42467, 28.53548, 28.90000,  
                        NA, 29.73818, 30.54667, 27.21652, 27.28800, 27.84000)
```

Instrução 1/4

- Obter a temperatura média do vetor `temperatura`

In []:

Instrução 2/4

- Obter as temperaturas que estão acima da média do vetor `temperatura`

In []:

Instrução 3/4

- Mostre quanto dos dados do vetor de temperaturas apresentam valores `NA`

In []:

Instrução 4/4

- Mostre onde os dados do vetor de temperaturas apresenta valores `NA`

In []:

9 Testes Lógicos com Vetores

- `any()` Testa se algum elemento do vetor atende a uma condição específica

Exemplo: Dado o vetor de nome `dap`, o qual armaneza dados de mensuração de diâmetro de milhares de árvores na Floresta Nacional de Altamira, teste se algum elemento é menor ou igual a 40.

```
In [ ]: any(dap >= 40)
```


- `all()` Testa se todos os elementos de um vetor atendem a uma condição.

Exemplo: Dado o vetor de nome `dap`, testar se algum elemento é menor do que 0:

```
In [ ]: all(dap < 40)
```

- `is.na()` Testa se o vetor contém valores ausentes (*Not Availables*)

```
In [ ]: vetor <- c(NA, 2, 3, 6)
is.na(vetor)
```

Exercício Prático - Índice de Vetor

Considerando o vetor de nome `dap`, o qual armaneza dados de mensuração de diâmetro de milhares de árvores na Floresta Nacional de Altamira, mostre:

Instrução 1/5

- Quantas árvores foram inventariadas.

```
In [ ]:
```

Instrução 2/5

- **Apenas** o penúltimo elemento desse vetor.

Instrução 3/5

- O diâmetro mínimo de medição

```
In [ ]:
```

Instrução 4/5

- O diâmetro máximo mensurado

```
In [ ]:
```

Instrução 5/5

- O diâmetro médio mensurado

```
In [ ]:
```