### Introdução ao R MATF14 - Estatística Econômica I

Rodney Fonseca

30/10/2024



#### Software R

R é uma linguagem e um ambiente para computação estatística e para preparação de gráficos

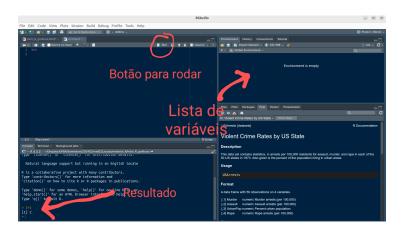
#### **Vantagens**

- software gratuito
- métodos para simulação probabilística e cálculo de probabilidades
- diversos tipos de técnicas para análise de dados
- ferramentas gráficas

#### Referências

- Livro Frery, Cribari-Neto (2011) Elementos de Estatística Computacional Usando Plataformas de Software Livre/Gratuito
- Página do curso de Ciência de dados da profa. Carolina Mota e prof. Gilberto Sassi do DEst-UFBA: https://ufba.netlify.app/paginas/catalogo
- Página do curso Estatística Computacional com R do prof. Paulo Justiniano (UFPR) e equipe: http://cursos.leg.ufpr.br/ecr/index.html

#### Ambiente do RStudio





## Operações básicas

#### Adição

```
1+1
```

```
## [1] 2
```

#### Subtração

```
5 - 3
```

```
## [1] 2
```

#### Multiplicação

```
2 * 3
```

```
## [1] 6
```

#### Divisão

### 10/5

## Operações básicas

```
Potência
2^4
## [1] 16
Logaritmo (natural)
log(10)
## [1] 2.302585
Exponencial
exp(3)
## [1] 20.08554
```

## Operadores lógicos

### Maior/menor que

```
5 > 3

## [1] TRUE

3>=3

## [1] TRUE

3>3

## [1] FALSE
```

### Operadores lógicos

#### **Igual**

```
5 == 3
## [1] FALSE
3 == 3
## [1] TRUE
Diferente
```

# 4 != 2

```
## [1] TRUE
```

```
4 != 4
```

```
## [1] FALSE
```

### Tipos de dados em R

```
Número real
class(0.5)
## [1] "numeric"
Valor lógico
class(TRUE)
## [1] "logical"
Caractere
class('UFBA')
## [1] "character"
```

### Variável

É como uma caixa nomeada. Você pode trocar o conteúdo da caixa, mas o nome permanece o mesmo.

Usamos os símbolos <- ou = para atribuir valores a uma variável

Para ver o valor da variável, basta rodar o seu nome

X

```
## [1] 2
```

#### Variável

Podemos fazer operações com a variável

```
5 * x

## [1] 10

podemos trocar o seu valor

x <- 20

x

## [1] 20
```

Podemos atribuir o valor de operações à outras variáveis

```
idade <- 2 * x
idade
```

```
## [1] 40
```

### Funções

 Em matemática, funções recebem um ou mais argumentos e retornam um ou mais valores. Funções em R são similares.

### Função para criar uma sequência de 9 números

```
seq(from = 1, to = 9)
## [1] 1 2 3 4 5 6 7 8 9
Fatorial de n
```

```
factorial(5)
```

```
## [1] 120
```

### Funções

 Em matemática, funções recebem um ou mais argumentos e retornam um ou mais valores. Funções em R são similares.

#### Função para criar uma sequência de 9 números

```
seq(from = 1, to = 9)
## [1] 1 2 3 4 5 6 7 8 9
```

#### Fatorial de n

```
factorial(5)
## [1] 120
```

#### Coeficiente binomial (combinação)

```
choose(5,2)
```

```
## [1] 10
```

Vetores são listas indexadas de variáveis do mesmo tipo. É como um armário de gavetas numeradas como 1, 2, 3, ... Você pode mudar o conteúdo da gaveta 5 sem alterar o conteúdo da gaveta 1.

Vetores podem ser criados com a fórmula c(x1, x2, ...)

```
meu_vetor \leftarrow c(1, 2, 5, 8, 1, 3)
```

Vetores são listas indexadas de variáveis do mesmo tipo. E como um armário de gavetas numeradas como 1, 2, 3, ... Você pode mudar o conteúdo da gaveta 5 sem alterar o conteúdo da gaveta 1.

Vetores podem ser criados com a fórmula c(x1, x2, ...)

```
meu_vetor \leftarrow c(1, 2, 5, 8, 1, 3)
```

Valores de elementos de um vetor podem ser vistos assim meu vetor[3]

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

Vetores são listas indexadas de variáveis do mesmo tipo. É como um armário de gavetas numeradas como 1, 2, 3, ... Você pode mudar o conteúdo da gaveta 5 sem alterar o conteúdo da gaveta 1.

```
Vetores podem ser criados com a fórmula c(x1, x2, ...)
```

```
meu_vetor \leftarrow c(1, 2, 5, 8, 1, 3)
```

Valores de elementos de um vetor podem ser vistos assim

```
meu_vetor[3]
```

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

Podemos alterar elementos do vetor

```
meu_vetor[3] <- 1
meu_vetor</pre>
```

```
## [1] 1 2 1 8 1 3
```

#### Soma dos elementos de um vetor

```
sum(meu_vetor)
## [1] 16
Produto dos elementos de um vetor
prod(meu vetor)
```

```
## [1] 48
```

#### Aplicando uma função aos elementos

```
log(meu_vetor)
```

```
[1] 0.0000000 0.6931472 0.0000000 2.0794415 0.0000000 1
```



### Funções para variáveis aleatórias

- R conta com funções para simulação e cálculo de probabilidades de diversas variáveis aleatórias: Bernoulli, binomial, geométrica, Poisson, normal, etc.
- Exemplos de funções: gerador de valores aleatórios, função de probabilidade, função de distribuição acumulada, etc.

### Funções para variáveis aleatórias

- R conta com funções para simulação e cálculo de probabilidades de diversas variáveis aleatórias: Bernoulli, binomial, geométrica, Poisson, normal, etc.
- Exemplos de funções: gerador de valores aleatórios, função de probabilidade, função de distribuição acumulada, etc.
- Vamos explorar algumas destas funções para as distribuições
   Bernoulli e binomial

- ▶ Relembrando: a distribuição Bernoulli é usada para modelar experimentos com dois resultados possíveis, *sucesso* ou *fracasso*, em que a probabilidade de sucesso é  $p \in [0,1]$
- Se  $X \sim Bernoulli(p)$ , então  $X \in \{0,1\}$  e X assume valor 1 com probabilidade p = P(X = 1).
- ▶ Temos ainda que E(X) = p e var(X) = p(1 p)

- ▶ Relembrando: a distribuição Bernoulli é usada para modelar experimentos com dois resultados possíveis, *sucesso* ou *fracasso*, em que a probabilidade de sucesso é  $p \in [0,1]$
- Se  $X \sim Bernoulli(p)$ , então  $X \in \{0,1\}$  e X assume valor 1 com probabilidade p = P(X = 1).
- ▶ Temos ainda que E(X) = p e var(X) = p(1 p)
- Note que também podemos dizer que  $X \sim Bin(1, p)$

Podemos usar a função  $rbinom(n\_amostras, 1, p)$  para gerar valores de  $X \sim Bernoulli(p)$ , em que  $n\_amostras$  é o número de realizações simuladas de X

- Podemos usar a função  $rbinom(n\_amostras, 1, p)$  para gerar valores de  $X \sim Bernoulli(p)$ , em que  $n\_amostras$  é o número de realizações simuladas de X
- **E**x.: Gerando 10 valores de  $X \sim Bernoulli(1/4)$ :

```
amostra_bern <- rbinom(10, 1, 1/4)
amostra_bern</pre>
```

```
## [1] 0 0 0 0 1 0 0 0 0
```

 Podemos montar tabelas para checar a frequência dos valores gerados

#### Tabela de frequências absolutas

```
table(amostra_bern)
```

```
## amostra_bern
## 0 1
## 9 1
```

#### Tabela de frequências relativas

```
table(amostra_bern)/10
```

```
## amostra_bern
## 0 1
## 0.9 0.1
```

Calculando a média amostral dos valores

```
mean(amostra_bern)
```

```
## [1] 0.1
```

Calculando a variância amostral dos valores

```
var(amostra_bern)
```

```
## [1] 0.1
```

Calculando a média amostral dos valores

```
mean(amostra_bern)
```

```
## [1] 0.1
```

Calculando a variância amostral dos valores

```
var(amostra_bern)
```

```
## [1] 0.1
```

**Observação:** as funções *mean* e *var* não fornecem os valores teóricos, somente estimativas com base na amostra gerada

Para calcular probabilidades, podemos usar a função **dbinom(k, 1, p)**. O exemplo abaixo fornece P(X = 1):

```
dbinom(1, 1, 1/4)
```

```
## [1] 0.25
```

Para calcular probabilidades, podemos usar a função **dbinom(k, 1, p)**. O exemplo abaixo fornece P(X = 1):

```
dbinom(1, 1, 1/4)
```

```
## [1] 0.25
```

A função de distribuição acumulada  $F(t) = P(X \le t)$  é **pbinom(t, 1, p)**. O exemplo abaixo fornece  $P(X \le 1/2)$ :

```
pbinom(1/2, 1, 1/4)
```

```
## [1] 0.75
```

▶ Vamos considerar um exemplo com 1000 valores simulados

```
amostra_maior_bern <- rbinom(1000, 1, 1/4)
table(amostra_maior_bern)/1000</pre>
```

```
## amostra_maior_bern
## 0 1
## 0.759 0.241
```

As frequências ficam mais próximas das probabilidades teóricas

Vamos considerar um exemplo com 1000 valores simulados

```
amostra_maior_bern <- rbinom(1000, 1, 1/4)
table(amostra_maior_bern)/1000</pre>
```

```
## amostra_maior_bern
## 0 1
## 0.759 0.241
```

- As frequências ficam mais próximas das probabilidades teóricas
- A média teórica é E(x) = p = 0.25, enquanto a amostral foi

```
mean(amostra_maior_bern)
```

```
## [1] 0.241
```

Quanto maior a amostra de valores gerados for, mais próxima a média amostral será da média teórica (populacional).

- Relembrando: a distribuição binomial conta o número de sucessos em n repetições de experimentos de Bernoulli com probabilidade p de sucesso
- ▶ Se  $X \sim Bin(n, p)$ , então  $X \in \{0, 1, ..., n\}$  e X assume valor  $k \in \{0, 1, ..., n\}$  com probabilidade

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

▶ Temos ainda que E(X) = np e var(X) = np(1-p)

▶ A função rbinom(n\_amostras, n, p) para gerar valores de X ~ Bin(n, p), em que n\_amostras é o número de realizações simuladas de X, enquanto n e p são os parâmetros da distribuição binomial

- ▶ A função rbinom(n\_amostras, n, p) para gerar valores de X ~ Bin(n, p), em que n\_amostras é o número de realizações simuladas de X, enquanto n e p são os parâmetros da distribuição binomial
- **E**x.: Gerando 10 valores de  $X \sim Bin(5, 1/2)$ :

```
amostra_bin <- rbinom(10, 5, 1/2)
amostra_bin</pre>
```

```
## [1] 3 1 2 3 3 2 4 3 3 4
```

Tabela de frequência dos valores gerados

```
table(amostra_bin)
```

```
## amostra_bin
## 1 2 3 4
## 1 2 5 2
```

- ▶  $X \sim Bin(n, p)$ , podemos calcular a probabilidade P(X = k) usando a função **dbinom(k, n, p)**.
- No exemplo abaixo, consideramos que  $X \sim Bin(5, 1/2)$  e calculamos P(X = 3):

```
dbinom(3, 5, 1/2)
```

```
## [1] 0.3125
```

▶ Usando a fórmula teórica  $P(X=3)=\binom{5}{3}(\frac{1}{2})^3(\frac{1}{2})^{5-3}$ , tal cálculo seria

```
choose(5, 3) * ( (1/2)^3 ) * ( (1/2)^(5 - 3) )
```

```
## [1] 0.3125
```

▶ Vamos gerar uma amostra grande de  $X \sim Bin(5, 1/2)$ :

```
tam_amostra <- 1000000
vx <- rbinom(tam_amostra, 5, 1/2)
table(vx)/tam_amostra</pre>
```

```
## vx
## 0 1 2 3 4 5
## 0.031465 0.156483 0.312334 0.312410 0.156142 0.031166
```

A média amostral dessa amostra é