# R para Ciência de Dados Intervalos de Confiança e Testes de Hipóteses

Profa Carolina Paraíba e Prof Gilberto Sassi

Departamento de Estatística Instituto de Matemática e Estatística Universidade Federal da Bahia

Junho de 2025

# Curso R para Ciência de Dados: Intervalos de Confiança e Testes de Hipóteses

Introdução

Revisão: Estatística Descritiva

Inferência Estatística

Revisão: Probabilidade

• Intervalos de Confiança

• Testes de Hipóteses

# Introdução

#### Durante o curso

- Usaremos nas aulas: posit.cloud.
- Recomendamos instalar e usar R com versão pelo menos 4.1: cran.r-project.org.
- usaremos o framework tidyverse:
  - Instalação: install.packages("tidyverse")

# Introdução

#### Na sua casa

- IDE recomendadas: RStudio e VSCode.
  - Caso você queira usar o VSCode, instale a extensão da linguagem R: REditorSupport.
- Outras linguagens interessantes: python e julia.
  - python: linguagem interpretada de próposito geral, contemporânea do R, simples e fácil de aprender.
  - julia: linguagem interpretada para análise de dados, lançada em 2012, promete simplicidade e velocidade.

#### Revisão de Estatística Descritiva no 'R'

#### Gráficos e Tabelas

#### Alguns conceitos básicos

- População: todos os elementos ou indivíduos alvo do estudo.
- Amostra: parte da população.
- Parâmetro: característica numérica da população. Usamos letras gregas para denotar parâmetros populacionais.
- Estatística: função ou cálculo da amostra
- Estimativa: característica numérica da amostra, obtida da estatística computada na amostra. Em geral, usamos uma estimativa para estimar o parâmetro populacional.
- Variável: característica mensurável comum a todos os elementos da população.

#### **Exemplo:**

- População: todos os eleitores nas eleições gerais de 2023.
- Amostra: 3.500 pessoas abordadas pelo datafolha.
- Variável: candidato a presidente de cada pessoa.
- Parâmetro: porcentagem de pessoas que escolhem Lula como presidente entre todos os eleitores.
- Estatística: porcentagem de pessoas que escolhem o lula
- **Estimativa:** porcentagem de pessoas que escolhem Lula como presidente entre todos os eleitores da amostra de 3.500 pessoas entrevistas pelo datafolha.

## Classificação de variáveis

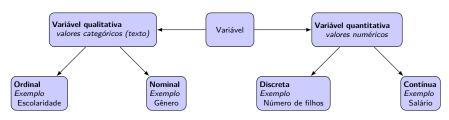


Figura 1: Classificação de variáveis.

# Tabela de distribuição de frequências Variável quantitativa discreta

A primeira coisa que fazemos é contar!

Χ	frequência	frequência relativa	porcentagem
$X_1$	$n_1$	$f_1$	100 · f <sub>1</sub> %
$X_2$	$n_2$	$f_2$	$100 \cdot f_2\%$
÷	:	:	÷
$X_k$	$n_k$	$f_k$	$100 \cdot f_k\%$
Total	n	1	100%

#### Observação

n é o tamanho da amostra.

## Tabela de distribuição de frequências

## Variável quantitativa discreta

```
dados <- read xlsx(</pre>
  "../dados/brutos/mulheres 20242.xlsx")
tab <- tabyl (dados, filhos) |>
  adorn totals() |>
  adorn pct formatting(digits = 2) |>
  rename (
    "Nro de filhos" = filhos,
    "Frequência" = n,
    "Porcentagem" = percent
```

#### tab

##

##

##

##

##	Nro	de	filhos	Frequência	Porcentagem
##			0	16	26.67%
##			1	7	11.67%
##			2	10	16.67%

6

12

60

3

4

5

Total

10.00%

20.00%

15.00%

100.00%

## Tabela de distribuição de frequências

# Variável quantitativa contínua

Para variáveis quantitativas discretas com muitos valores distintos, e para variáveis quantitativas contínuas.

X	frequência	frequência relativa	porcentagem
$[l_0, l_1)$	$n_1$	$f_1$	$100 \cdot f_1\%$
$[I_1, I_2)$	$n_2$	$f_2$	$100 \cdot f_2\%$
$[I_2, I_3)$	$n_3$	$f_3$	$100 \cdot f_3\%$
:	:	:	:
$[I_{k-1},I_k]$	$n_k$	$f_k$	$100 \cdot f_k\%$
Total	n	1	100%

#### Observação

n é o tamanho da amostra.

#### Tabela de distribuição de frequências

# Variável quantitativa contínua

```
dados <- read xlsx("../dados/brutos/iris.xlsx")</pre>
dados <- clean names (dados)
k <- floor(1 + log2(nrow(dados)))</pre>
dados <- dados |>
  mutate(comprimento_sepala_int = cut(
    comprimento_sepala,
    breaks = k,
    include.lowest = TRUE,
    right = FALSE
  ) )
```

```
tab <- tabyl(dados, comprimento_sepala_int) |>
  adorn_totals() |>
  adorn_pct_formatting(digits = 2) |>
  rename(
    "Comprimento de Sépala" = comprimento_sepala_int,
    "Frequência" = n,
    "Porcentagem" = percent
)
```

# tab

##

##

##	Comprimento	ae	Sepala	Frequencia	Porcentager
##		[4.3	3,4.75)	11	7.339
##		[4.	75,5.2)	30	20.009
##		[5.2	2,5.65)	24	16.009
##		[5.6	65,6.1)	24	16.009
##		[6.3	1,6.55)	31	20.679
##		[ (	6.55,7)	17	11.339
##		٦٦	7,7,45)	7	4.679

Total

4.00%

100.00%

150

[7.45,7.9]

## Histograma

Para variávieis quantitativas contínuas, geralmente não construímos gráficos de barras, e sim uma figura geométrica chamada de *histograma*.

- O histograma é um gráfico de barras contíguas em que a área de cada barra é igual à frequência relativa.
- Cada faixa de valor  $[l_{i-1}, l_i)$ , i = 1, ..., n, será representada por um barra com área  $f_i$ , i = 1, ..., n.
- Como cada barra terá área igual a  $f_i$  e base  $l_i l_{i-1}$ , e a altura de cada barra será  $\frac{f_i}{l_i l_{i-1}}$ .
- $\frac{f_i}{l_i-l_{i-1}}$  é denominada de densidade de frequência.

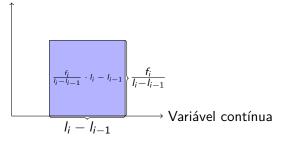
## Histograma

- Podemos usar os seguintes parâmetros (obrigatório o uso de apenas um deles):
  - bins: número de intervalos no histograma (usando, por exemplo, a regra de Sturges).
  - binwidth: tamanho (ou largura) dos intervalos.
  - breaks: os limites de cada intervalo.

#### Histograma

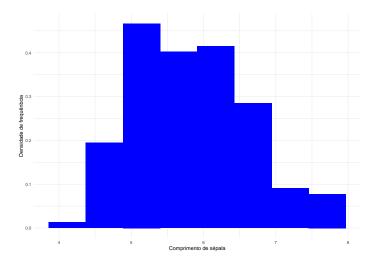
Figura 2: Representação de uma única barra de um histograma.

Denside de frequência



#### Histograma

```
ggplot(dados) +
  geom_histogram(
    aes(comprimento_sepala, after_stat(density)),
    bins = k,
    fill = "blue"
) +
labs(
    x = "Comprimento de sépala",
    y = "Densidade de frequênbcia"
) +
theme_minimal()
```



#### Medidas Resumo

```
tab <- group_by(dados, especies) |>
summarise(
   media = mean(comprimento_sepala),
   dp = sd(comprimento_sepala),
   cv = dp / media,
   q1 = quantile(comprimento_sepala, probs = 1 / 4),
   q2 = quantile(comprimento_sepala, probs = 2 / 4),
   q3 = quantile(comprimento_sepala, probs = 3 / 4)
)
```

#### tab

3 virginica 6.59 0.636 0.0965 6.22 6.5

6.9

#### Inferência Estatística

## O que faremos nesse curso?

• Estimação Pontual: Aproximar um parâmetro.

Exemplo: Estimar o teor alcóolico de uma bebida.

 Intervalo de Confiança: Encontrar uma estimativa intervalar para um parâmetro.

Exemplo: Encontrar números a e b tal que o teor alcóolico verdadeiro está entre a e b com uma probabilidade estabelecida pelo pesquisador.

#### Inferência Estatística

#### O que faremos nesse curso?

• **Teste de Hipóteses:** Decidir entre duas hipóteses  $H_0$  e  $H_1$ : negação de  $H_0$ .

Exemplo: Decidir entre duas hipóteses:

 ${\it H}_{\rm 0}$  : A nota média em matemática no ENEM 2021 é maior que 600,

 $H_1$ : A nota média em matemática no ENEM 2021 é menor ou igual 600.

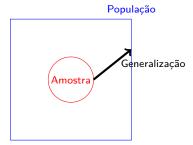
#### Inferência Estatística

Em todos estes casos, precisamos usar *Probabilidade*.

#### Por que precisamos de probabilidade?

- Queremos fazer afirmações válidas para toda população.
- Inferência Estatística: generalização da amostra para toda população precisa de probabilidade.

Figura 3: Ilustração da estatística inferencial.



#### **Probabilidade**

#### Fenômeno Aleatório

Procedimento ou evento cujo resultado não é possível antecipar de forma determinística. Por exemplo:

• Qual será o resultado do lançamento de um dado "justo"?

#### Espaço amostral

O conjunto de todos os resultados de um fenômeno aleatório. Notação:  $\Omega$ .

#### **Evento**

Subconjunto de um espaço amostral. Notação: A, B, C, . . . .

#### Ponto amostral

Um resulto possível de um fenómeno aleatório. Notação:  $\omega$ .

#### **Probabilidade**

#### **Probabilidade**

A plausibilidade de um ponto amostral  $\omega$  de A ser o resultado do fenômenos aleatório. Notação: P(A).

#### Variável aleatória

Função com domínio em um espaço amostra e contra-domínio no conjunto dos números reais  $X:\Omega\to\mathbb{R}.$ 

#### **Probabilidade**

#### Classificação de variáveis aleatórias

- Dizemos que X é uma variável aleatória discreta, se os valores possíveis desta variável são números inteiros, geralmente resultado de contagem;
- Dizemos que X é uma variável aleatória contínua, se os valores possíveis desta variável pode ser qualquer número (incluindo aqueles por parte decimal);
- O conjunto dos valores possíveis de X representamos por  $\chi$ .

#### Variável Aleatória Discreta

#### Função de probabilidade (FP):

$$f(x) = P(X = x)$$

**Interpretação:** f(x) pode ser interpretada como a frequência relativa de x em toda população.

#### Na amostra

X	frequência relativa
<i>x</i> <sub>1</sub>	$f_1$
$x_2$	$f_2$
<i>X</i> 3	$f_3$
:	:
$x_k$	$f_k$

## Na população

X	função de probabilidade
$\overline{x_1}$	$f(x_1)$
<i>x</i> <sub>2</sub>	$f(x_2)$
<i>X</i> 3	$f(x_3)$
:	<u>:</u>
$X_k$	$f(x_k)$

## Variável Aleatória Discreta

## Medidas resumo para variável aleatória discreta

#### Na amostra

X uma variável quantitativa discreta

Média:

$$\bar{X} = x_1 \cdot f_1 + \cdots + x_k \cdot f_k$$

Variância:

$$Var(X) = (x_1 - \bar{x})^2 \cdot f_1 + \dots + (x_k - \bar{x})^2 \cdot f_k$$

• Desvio padrão:

$$dp(X) = \sqrt{Var(X)}$$

• Mediana:

Md tal que:

- $f_1 + \cdots + f_{Md} \ge 0,5$
- $f_{Md} + \cdots + f_k \leq 0, 5$

#### Na população

X para uma variável aleatória discreta

Média:

$$\mu = x_1 \cdot f(x_1) + \cdots + x_k \cdot f(x_k)$$

• Variância:

$$\sigma^2 = (x_1 - \mu)^2 \cdot f(x_1) + \dots + (x_k - \mu)^2 \cdot f(x_k)$$

• Desvio padrão:

$$\sigma = \sqrt{Var(X)}$$

Mediana:

Md tal que:

• 
$$f(x_1) + \cdots + f(Md) \ge 0,5$$

• 
$$f(Md) + \cdots + f(x_k) \le 0,5$$

# Distribuição Bernoulli

Cada elemento da população pode ter sucesso ou fracasso.

**Sucesso:** caso de interesse ou mais importante.

Sucesso	Fracasso	
Munícipio tem secretaria cultura	Munícipio <b>não</b> tem secretaria cultura	
Pessoa infectada	Pessoa sadia	
Pessoa alta	Pessoa baixa	
Bahia ganha o jogo	Bahia <b>não</b> ganha o jogo	

Precisamos descobrir a proporção (ou porcentagem) de Sucesso.

Notação: p é a prporção (ou porcentagem) de Sucesso.

# Distribuição Bernoulli

#### Parâmetros da distribuição Bernoulli

Usamos letras gregas para representar parâmetros:

- Média populacional:  $\mu$
- Variância populacional:  $\sigma^2$
- Desvio padrão populacional:  $\sigma$

#### Distribuição Bernoulli

- Média (populacional):  $\mu = p$
- Variância (populacional):  $\sigma^2 = p \cdot (1-p)$
- Desvio padrão (populacional):  $\sigma = \sqrt{p \cdot (1-p)}$

# Estimação pontual Disribuição Bernoulli

- Definimos o sucesso.
- 2 Encontramos a estimativa de p.

Variável aleatória: transmissão (do conjunto de dados mtcarros.xlsx).

- 0: Carro com transmissão automática
- 1: Carro com transmissão manual (Sucesso)

# Estimação pontual Disribuição Bernoulli

# Estimação pontual Disribuição Bernoulli

Variável aleatória: Mulherer já teve gravidez? (coluna gravidez em mulheres\_20242.xlsx).

- Sucesso: Sim (sim a mulher já esteve grávida).
- Fracasso: Não (nao a mulhere nunca esteve grávida).

Vamos criar uma nova coluna com 1 e 0.

```
dados <- read xlsx(
  "../dados/brutos/mulheres 20242.xlsx")
dados <- dados |>
 mutate(ind_gravidez = ifelse(gravidez == "sim", 1, 0))
tab <- dados |>
  summarise(prop_gravidez = mean(ind_gravidez))
tab
## # A tibble: 1 x 1
## prop_gravidez
```

## <dbl>

# Distribuição Binomial

- Temos n casos
- Cada caso pode ser sucesso ou fracasso

#### Parâmetros:

- Proporção de sucesso: p
- Número de casos: n
- Média:  $\mu = n \cdot p$
- Variância:  $\sigma^2 = n \cdot p \cdot (1 p)$
- Desvio padrão:  $\sigma = \sqrt{n \cdot p \cdot (1-p)}$

Precisamos estimar p.

Geralmente conhecemos previamente n.

Soma de Bernoulli produz Binomial.

# Estimação pontual Distribuição Binomial

Variável aleatória: Número mulheres que com diagnóstico de endometriose (coluna endometriose de mulheres\_20242.xlsx).

```
dados <- read_xlsx(
   "../dados/brutos/mulheres_20242.xlsx")

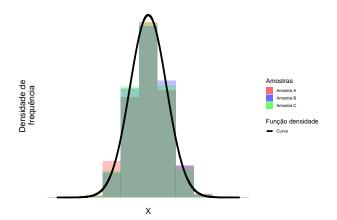
tab <- dados |>
   summarise(
    prop = sum(endometriose) / nrow(dados))
tab
```

```
## # A tibble: 1 x 1
## prop
## <dbl>
## 1 0.167
```

#### Variável Aleatória Contínua

#### Motivação

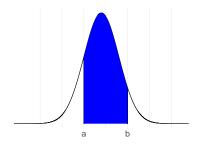
- Para cada amostra, temos um histograma;
- Queremos encontrar uma curva que aproxima bem todos os histogramas possíveis;



#### Variável Aleatória Contínua

#### Propriedades de variável aleatória contínua

- Proporção de elementos da população com variável aleatória X entrea a e b: P(a < X < b).
- P(a < X < b): área sob a curva (região azul).



### Distribuição Normal

#### Quando usar?

- Valores da variável aleatória concentrados em torno da média;
- Valores da variável aleatória afastados da média são pouco prováveis;
- Função densidade de probabilidad em curva em formato de sino;
- Simetria em torno da média.

#### Checamos isso com histograma.

#### **Parâmetros**

Média: μ

• Variância:  $\sigma^2$ 

• Desvio padrão:  $\sigma$ 

## Exemplos de aplicação Distribuição normal

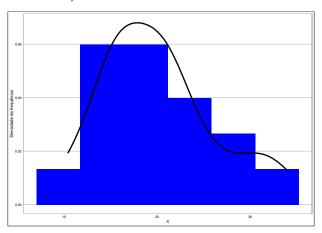
- Altura:
  - As pessoas no Brasil tem em média 170cm
  - Algumas pessoas são menores que 170cm
  - Algumas pessoas são maiores que 170cm
  - poucas ficam muito longe de 170cm
- Uso de caixa eletrônico:
  - Em média, as pessoas demoram 2 minutos no caixa eletrônico
  - Algumas pessoas são mais lentas
  - Algumas pessoas são mais rápidas
  - poucas pessoas ficam longe de 2 minutos

# Exemplos Distribuição normal

Variável aleatória: milhas por galão (milhas\_por\_galao em mtcarros.xlsx).

# Exemplos Distribuição normal

Variável aleatória: milhas por galão (milhas\_por\_galao em mtcarros.xlsx).



## Estimativa pontual Distribuição Normal

## media dp ## <dbl> <dbl> ## 1 20.1 6.03

## Intervalo de Confiança

#### Intervalo de Confiança

**Objetivo:** Para parâmetro  $\mu$  ( $\sigma$  e p), encontrar L e U tal que  $L < \mu < U$  com alguma probabilidade associada  $\gamma$ .

Chamamos  $\gamma$  de coeficiente de confiança.

Vamos usar o pacote statBasics.

#### Interpretação de Intervalo de Confiança

O parâmetro  $\mu$  ( $\sigma$  e p) pode ou não pode estar entre L e U do intervalo de confiança com coeficiente de confiança  $\gamma$ .

```
dados_estudos <- read_xlsx("../dados/brutos/motivacao_intervalo_confianca.xlsx", sheet = 1)
dados_pop <- read_xlsx("../dados/brutos/motivacao_intervalo_confianca.xlsx", sheet = 2)
media_pop <- mean(dados_pop$variavel)
tab <- dados estudos |>
```

group by (estudo) |>

media pop = media pop)

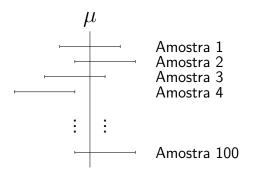
#### tab

```
## # A tibble: 6 x 4
## estudo lower_ci upper_ci media_pop
## <chr>
            <dbl>
                    <dbl>
                            <dbl>
                            6.37
## 1 amostra1
              5.40
                 6.30
                           6.37
## 2 amostra2
             4.53 6.17
## 3 amostra3
                           6.37
             5.76 6.19
## 4 amostra4 5.74 8.29
                           6.37
## 5 amostra5
             5.63 7.99
                           6.37
## 6 amostra6 4.12 8.70
                           6.37
```

## Intervalo de Confiança

 $\gamma\%$  das amostras vão gerar intervalos de confiança que contém o parâmetro.

Figura 4: Interpretação de intervalo de confiança.



Geralmente  $\gamma$  é 99%, 95% ou 90%.

# Intervalo de Confiança Distribuição Bernoulli

**Primeira forma:** Vetor de 1 e 0.

- Variável aleatória: Carro tem transmissão manual? (variável transmissao de mtcarros.xlsx).
- Sucesso: 1 (transmissão manual)
- Fracasso: 0 (transmissão automática)

A prorporção de carros com transmissão manual está entre 0,233 e 0,5795 com coeficiente de confiança 95%.

# Intervalo de Confiança Distribuição Bernoulli

**Segunda forma:** Número de tentativas e número de sucessos.

- Variável aleatória: Carro tem transmissão manual? (variável transmissao de mtcarros.xlsx).
- Sucesso: 1 (transmissão manual)
- Fracasso: 0 (transmissão automática)

A proporção de carros com transmissão automática está entre 0,233 e 0,5795 com coeficiente de confiança 95%.

# Intervalo de Confiança Distribuição Binomial

Pesquisa de Intenção de voto: Eleição 2023.

Número de entrevistados: 8308Número de eleitores de Lula: 4403

Coeficiente de Confiança: 99%

```
eleicao_lula_22 <- ci_lpop_bern(4403, 8308, conf_level = 0.99)
eleicao_lula_22

## # A tibble: 1 x 3
## lower_ci upper_ci conf_level
## <dbl> <dbl> <dbl>
## 1 0.516 0.544 0.99
```

Lula teria uma proporção entre 0,5158 e 0,5441 de votos com coeficiente de 99%.

# Intervalo de Confiança Distribuição Binomial

- Variável aleatória: Número de mulheres que já engravidaram (coluna gravidez de mulheres\_20242.xlsx);
- Coeficiente de confiança: 92,5%.

A proporção de mulheres que já engravidaram está entre 0,6184 e 0,8483 com coeficiente de confiança 92,5%.

## Intervalo de Confiança para média

- Variável aleatória tem distribuição normal;
- Variância populacional desconhecida;
- Intervalo de confiança para média.

Função ci\_1pop\_norm do pacote statBasics.

### Intervalo de Confiança para média

- Variável aleatória: milhas por galão (milhas\_por\_galao em mtcarros.xlsx).
- Coeficiente de confiança: 99%.

Os carros fazem, em média, entre 17,17 e 23,01 milhas por galão com coeficiente de confiança 99%.

## Intervalo de Confiança para variância

- Variável aleatória tem distribuição normal;
- Média populacional é desconhecida;
- Intervalo de confiança para a variância.

Função ci\_1pop\_norm do pacote statBasics, com parâmetro parameter='variance'.

## Intervalo de Confiança para variância

- Variável aleatória: milhas por galão (milhas\_por\_galao em mtcarros.xlsx).
- Coeficiente de confiança: 99%.

A variabilidade do consumo de milhas por galão está entre 20,47 e 77,89, com coeficiente de confiança 99%.

## Intervalo de Confiança para média Grandes amostras

- Variável aleatória não tem distribuição normal;
- · Variância populacional é desconhecida;
- Tamanho amostral é suficientemente grande;
- Intervalo de confiança para média.

Função ci\_1pop\_general do pacote statBasics.

## Intervalo de Confiança para média Grandes amostras

- Variável aleatória: número de filhos (filhos em mulheres\_20242.xlsx).
- Coeficiente de confiança: 95%.

As mulheres têm, em média, entre 1,824 e 2,776 filhos com coeficiente de confiança 95%.

#### Objetivo:

Decidir entre  $H_0$  (hipótese nula) e  $H_1$  (hipótese alternativa) usando as evidências da amostra.

- H<sub>0</sub> é a negação de H<sub>1</sub>
- $H_1$  é a negação de  $H_0$
- *H*<sub>1</sub> é aquilo que desejamos provar que é verdade
  - H<sub>1</sub> é afirmação extraordinária que precisa de evidências para acreditarmos
- H<sub>0</sub> é o padrão, valor padrão de mercado ou valor padrão do regulador (ex. ANVISA)
  - H<sub>0</sub> é a afirmação ordinária que assumimos como verdade quando não temos evidência para acreditar em H<sub>1</sub>

- Decisão através de evidência na amostra:
  - Decisão baseada em  $\it evidência \implies hipótese alternativa <math>\it H_1$
  - Decisão  $sem\ evidência\$ ou  $na\ dúvida \Longrightarrow$ hipótese nula  $H_0$

Como temos uma **tendência de continuar em**  $H_0$  na ausência de *evidências*, escrevemos:

- Decisão por H<sub>0</sub>: Não rejeitamos H<sub>0</sub>;
- Decisão por H<sub>1</sub>: Não rejeitamos H<sub>1</sub>.

#### Podemos cometer dois erros ao decidir:

- Erro tipo I ou Falso positivo: Decisão por H<sub>1</sub>, mas H<sub>0</sub> é a verdade. Erro GRAVÍSSIMO!
- Erro tipo II ou Falso negativo: Decisão por H<sub>0</sub>, mas H<sub>1</sub> é a verdade
- Nível de significância:  $\alpha = P(Falso positivo)$
- **Poder do teste**:  $1 \beta = P(Verdadeiro positivo)$

		Situação na população	
		H <sub>0</sub>	$H_1$ (Negação de $H_0$ )
Decisão	H <sub>0</sub>   H <sub>1</sub> (Negação de H <sub>0</sub> )	Sem erro (verdadeiro negativo) Falso positivo (Erro tipo I)	Falso negativo (Erro tipo II) Sem erro (Verdadeiro positivo)

**Objetivo:** Como  $H_1$  (positivo) é a hipótese mais importante, então queremos decidir entre  $H_0$  e  $H_1$  garantindo que:

- o *nível de significância* seja pequeno (geralmente 5%)
- o poder do teste seja máximo possível

- Sem evidência, continuamos acredidanto em  $H_0$ .
- Com evidência, desistimos de  $H_0$  e passamos a acreditar em  $H_1$ .

Neste contexto, usamos o verbo rejeitar em estatística:

- Sem evidência, não rejeitamos H<sub>0</sub>.
- Com evidência, rejeitamos H<sub>0</sub>.

# Teste de Hipóteses Exemplo

Em um julgamento, temos as seguintes hipóteses:

H<sub>0</sub>: o réu é inocente
H<sub>1</sub>: o réu é culpado

Em um julgamento, o sistema de justiça pode cometer dois erros:

- falso positivo: uma pessoa inocente é condenada
- falso negativo: um pessoa culpada é inocentada

Em um julgamento, o sistema de justiça usa a seguinte regra de decisão:

- réu é culpado: apenas se tiver evidências fortes e concretas
- réu é inocente: na dúvida ou na ausência de evidências

# Teste de Hipóteses Como decidir?

Hipóteses nula e alternativa geralmente são declarações matemáticas envolvendo parâmetros.

**Ideia:** Calculamos uma distância entre a estimativa e o valor do parâmetro quando a hipótese nula é verdade.

- Se essa distância for pequena, decidimos por  $H_0$
- Se essa distância for grande, decidimos por  $H_1$

Chamamos esta distância de Estatística Teste.

Existem duas formas de determinar o que é pequeno ou grande (extrema):

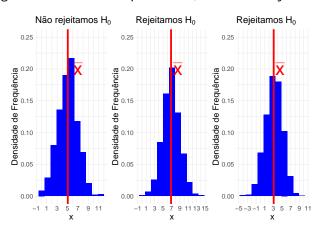
- 1 Procedimento Geral de Testes de Hipóteses ou Procedimento de Neymann-Pearson
- **2** valor-p (*p-value* em inglês)

## Teste de Hipóteses Como decidir?

População:  $N(\mu, 4)$ .

• Hipóteses:  $H_0$ :  $\mu = 5$  contra  $H_1$ :  $\mu \neq 5$ .

• Regra de Decisão: se  $\bar{x}$  perto de 5, então não rejeitamos  $H_0$ .



## Procedimento de Neymann-Pearson

#### **Etapas:**

- **1** Estabeleça  $H_0$  e  $H_1$
- 2 Esteleça o (máximo) nível de significância
- 3 Encontre a *região crítica* (conjunto onde a *estatística teste* é grande)
- 4 Verifique se a estatística teste está na região crítica

A região crítica é construída usando o nível de significância.

Erros decaem quando o tamanho amostral aumenta.

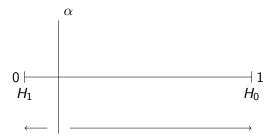
As porcentagens de falso positivo e de falso negativo diminuem quando o tamanho da amostra aumenta.

## valor-p ou nível crítico p-value

- Valor-p é uma medida de evidência contra o hipótese nula.
- Valor-p NÃO É A PROBABILIDADE DO FALSO POSITIVO.
- Para cada amostra, temos um valor-p diferente.
- Formalmente: probabilidade de coletar uma outra amostra (de mesmo tamanho) com estatística de teste mais extrema do que a amostra que eu tenho se a hipótese nula é verdadeira.
- Rejeitamos H<sub>0</sub> se o valor-p for menor que o nível de significância.

# valor-p ou nível crítico p-value

Rejeitamos o valor-p menor que nível de significância:  $p < \alpha$ .



# valor-p ou nível crítico p-value

#### **Etapas:**

- ① Estabeleça  $H_0$  e  $H_1$
- 2 Esteleça o (máximo) nível de significância
- 3 Calcule o valor-p
- 4 Verifique se o valor-p é menor que nível de significância

Para cada amostra, temos um valor-p diferente.

O valor-p (p) pode ser pequeno ou pode ser grande.

Se  $H_0$  é verdade, aproximadamente  $\alpha$ % das amostras produzem o falso positivo quando usamos o valor-p.

### Teste t para média

- Variável aleatória tem distribuição normal;
- Variância populacional é desconhecida;
- Hipóteses sobre a média da população  $(\mu)$ .

No pacote statBasics: ht\_1pop\_mean.

#### Testes de hipóteses deste curso:

- Teste unilateral à esquerda:  $H_1$ :  $\mu < \mu_0$ 
  - alternative = 'less'
- Teste unilateral à direita:  $H_1$ :  $\mu > \mu_0$ 
  - alternative = 'greater'
- Teste bilateral:  $H_1$ :  $\mu \neq \mu_0$ 
  - alternative = 'two.sided' valor padrão

Especificamos  $\mu_0$  com o parâmetro mu.

### Teste t para média

Temos evidência para afirmar que os carros americanos conseguem fazer no máximo 15 milhas por galão, em média, ao nível de significância 1%?

- $H_0$  (negação de  $H_1$ ):  $\mu \ge 15$
- $H_1$  (o que queremos provar):  $\mu < 15$

## Teste t para média

```
dados <- read_xlsx(".../dados/brutos/mtcarros.xlsx")
ht_milhas_galao <- ht_lpop_mean(
   dados$milhas_por_galao,
   mu = 15,
   alternative = "less",
   sig_level = 0.01)
ht_milhas_galao

## # A tibble: 1 x 7
## statistic p value critical value critical region alternative mu sig level</pre>
```

Não há evidência para afirmar que os carros americanos fazer no máximo 15 milhas por galão, em média, ao nível de significância 1%.

- Hipóteses sobre a proporção de sucessos (p);
- Variável aleatória tem distribuição Bernoulli ou distribuição binomial.

#### Testes de hipóteses deste curso:

- Teste unilateral à esquerda:  $H_1$ :  $p < p_0$ 
  - alternative = 'less'
- Teste unilateral à direita:  $H_1$ :  $p > p_0$ 
  - alternative = 'greater'
- Teste bilateral:  $H_1: p \neq p_0$ 
  - alternative = 'two.sided' valor padrão

Especificamos  $p_0$  com o parâmetro proportion.

## Teste z para proporção Distribuição Bernoulli

Temos evidência para afirmar que a proporção de carros americanos com transmissão manual é maior que 25% ao nível de significância 1%?

- $H_0$  (negação de  $H_1$ ):  $p \le 0,25$
- $H_1$  (o que desejamos provar): p > 0,25

## Teste z para proporção Distribuição Bernoulli

```
dados <- read xlsx("../dados/brutos/mtcarros.xlsx")</pre>
teste transmissao <- ht 1pop prop (
  dados$transmissao,
  proportion = 0.25,
  alternative = 'greater',
  sig level = 0.01)
teste transmissao
## # A tibble: 1 x 7
## statistic p_value critical_value critical_region alternative proportion
    <dbl> <dbl>
                       <dhl> <chr>
                                         <chr>
                                                      <dh1>
## 1 2.04 0.0206
                       2.33 (2.326, Inf) greater
                                                      0.25
## # i 1 more variable: sig level <dbl>
```

Ao nível de significância 1%, não temos evidência para afirmar que a proporção de carros americanos com transmissão automática é maior 25%.

## Teste z para proporção Distribuição Binomial

- Variável aleatória: Número de mulheres com diagnóstoco de endometriose (coluna endometriose de mulheres\_20242.xlsx)
- Nível de significância: 2,5%
- A maioria das mulheres têm endometriose?
  - $H_0: p < 0.5$
  - $H_1: p > 0,5$  (o que desejamos provar)

## Teste z para proporção Distribuição Binomial

```
dados <- read_xlsx("../dados/brutos/mulheres_20242.xlsx"
n_tentativas <- nrow(dados)
n_sucessos <- sum(dados$endometriose)
teste_endometriose <- ht_lpop_prop(
    n_sucessos, n_tentativas, proportion = 0.5,
    alternative = "greater", sig_level = 0.025)

teste_endometriose

## # A tibble: 1 x 7
## statistic p_value critical_value critical_region alternative proportion
## <dbl> <dbl > dbl> <dbl> <dbl > dbl> <dbl > dbl> <dbl > dbl> <dbl > dbl > dbl> <dbl > dbl > dbl
```

Ao nível de significância 2,5%, a maior parte das mulheres têm endometriose.

# Teste z para média Grandes amostras

- Hipóteses sobre a média da população  $(\mu)$ ;
- Variável aleatória não tem distribuição normal.
- Tamanho amostral é suficientemente grande;

No pacote statBasics: ht\_1pop\_mean (com adaptação).

## Teste Qui-Quadrado para variância Distribuição normal

- Hipóteses sobre a variância da população  $(\sigma)$ ;
- Variável aleatória tem distribuição normal.
- Média populacional é desconhecida.

No pacote statBasics: ht\_1pop\_var.

#### Testes de hipóteses deste curso:

- Teste unilateral à esquerda:  $H_1$ :  $\sigma^2 < \sigma_0^2$ 
  - alternative = 'less'
- Teste unilateral à direita:  $H_1$ :  $\sigma^2 > \sigma_0^2$ 
  - alternative = 'greater'
- Teste bilateral:  $H_1: \sigma^2 \neq \sigma_0^2$ 
  - alternative = 'two.sided' valor padrão

Especificamos  $\sigma_0^2$  com o parâmetro sigma.

## Teste Qui-Quadrado para variância Distribuição normal

Temos evidência para afirmar o desvio padrão da distância percorrida por galão nos carros americanos é menor 2 (milhas por galão), ao nível de significância 5%?

•  $H_0$  (negação de  $H_1$ ):  $\sigma^2 \ge 2^2$ 

```
• H_1 (o que queremos provar): \sigma^2 < 2^2 dados <- read_xlsx("../dados/brutos/mtcarros.xlsx") ht_milhas_galao <- ht_lpop_var( dados$milhas_por_galao, alternative = "less", sigma = 2, sig_level = 0.05)
```

```
ht_milhas_galao
```

```
## # A tibble: 1 x 7
## statistic p_value critical_value critical_region alternative sigma sig_level
## <dbl> <dbl> <dbl> <chr> <dbl> <chr> <dbl> = dbl> <dbl> <chr> <dbl> = dbl> <dbl> <chr> <dbl> = dbl> <dbl> <d
```

Ao nível de significância 5%, não temos evidência para afirmar o desvio padrão da distância percorrida por galão nos carros americanos é menor 2 (milhas por galão).

#### **Experimentos Comparativos**

#### Experimento completamente aleatório:

Medimos uma mesma variável em duas populações independentes.

- 1 População 1
- 2 População 2
- 3 As duas populações são independentes

#### Estudo observacional:

- 1 Acompanhamos cada elemento da amostra *antes* e *depois* de uma *intervenção*
- 2 As duas populações não são independentes
- Teste t pareado

### Comparação de variâncias

Antes de comparar  $\mu_1$  e  $\mu_2$ , precisamos verificar se  $\sigma_1=\sigma_2$ 

- População 1:  $N(\mu_1, \sigma_1^2)$
- População 2:  $N(\mu_2, \sigma_2^2)$
- Teste de Hipóteses envolvendo  $\sigma_1$  e  $\sigma_2$

No pacote statBasics: ht\_2pop\_var.

### Comparação de variâncias

Testes de hipóteses deste curso:

- Teste bilateral:  $H_1: \sigma_1 \neq \sigma_2$ 
  - alternative ='two.sided'- valor padrão
- Teste unilateral à esquerda:  $H_1: \sigma_1 < \sigma_2$ 
  - alternative = 'less'(Atenção para ordem das populações)
- Teste unilateral à direita:  $H_1: \sigma_1 > \sigma_2$ 
  - alternative = 'greater' (Atenção para ordem das populações)

Especificamos ratio fornecendo  $\frac{\sigma_1}{\sigma_2}$ . Valor padrão: ratio = 1 (neste caso, estamos testando a igualdade).

### Comparação de variâncias

Ao nível de signficância 1%, existe diferença entre os desvios padrões da distância percorrida em milhas por um galão entre carros com transmissão manual e automática.

- Variável aleatória: Milhas por galão
- População 1: carros com transmissão manual (transmissão == 1)
- População 2: carros com transmissão manual (transmissão == 0)

#### Comparação de variâncias

```
dados <- read_xlsx("../dados/brutos/mtcarros.xlsx")
carros_manuais <- dados |> filter(transmissao == 1)
carros_auto <- dados |> filter(transmissao == 0)
comparacao_var <- ht_2pop_var(
   carros_manuais$milhas_por_galao,
   carros_auto$milhas_por_galao,
   ratio = 1, sig_level = 0.01)
comparacao_var</pre>
```

```
## # A tibble: 2 x 7

## statistic p_value critical_vale ratio alternative lower_ci upper_ci
## <dbl> <dbl> <dbl> <dbl> <chr> <dbl> <dbl> <chr> <dbl> 2.59 0.0669 0.218 1 two.sided 2.64 2.64
## 2 2.59 0.0669 3.86 1 two.sided 2.64 2.64
```

Não temos evidência para assumir que as variâncias são diferentes ao nível de significâcia 1%, e assumimos que as variâncias são iguais.

Primeiro precisamos verificar se os desvios padrões são iguais para duas populações.

- Variável aleatória: milhas percorridas por galão (milhas\_por\_galao)
- População 1: carros com transmissão manual (transmissão == 1)
- População 2: carros com transmissão automática (transmissão == 0)

```
carros_manuais <- dados |> filter(transmissao == 1)
carros_auto <- dados |> filter(transmissao == 0)
comparacao_var <- ht_2pop_var(
    carros_manuais$milhas_por_galao,
    carros_auto$milhas_por_galao)
comparacao_var

## # A tibble: 2 x 7</pre>
```

dados <- read xlsx(".../dados/brutos/mtcarros.xlsx")</pre>

Ao nível de significância 5%, continuamos acreditando que os desvios padrões das duas populações são iguais.

Quando sabemos que as variâncias populacionais são iguais.

- População 1: N(μ<sub>1</sub>, σ)
- População 1:  $N(\mu_2, \sigma)$
- Teste de Hipóteses envolvendo  $\mu_1$  e  $\mu_2$

No pacote statBasics: ht\_2pop\_mean com argumento var\_equal = T.

#### Testes de hipóteses deste curso:

- Teste bilateral:  $H_1: \mu_1 \mu_2 = \Delta_0$ 
  - alternative = 'two.sided'- valor padrão
- Teste unilateral à esquerda:  $H_1: \mu_1 \mu_2 < \Delta_0$ 
  - alternative = 'less' (Atenção para ordem das populações)
- Teste unilateral à direita:  $H_1: \mu_1 \mu_2 > \Delta_0$ 
  - alternative = 'greater' (Atenção para ordem das populações)

Especificamos delta fornecendo  $\Delta_0 = \mu_1 - \mu_2$ . Valor padrão: delta = 0 (neste caso, estamos testando a igualdade).

0 greater

Ao nível de significância 1%, carros com transmissão automática andam mais com galão de gasolina que carros com transmissão manual?

```
comparacao_medias <- ht_2pop_mean(
  carros_auto$milhas_por_galao,
  carros_manuais$milhas_por_galao,
  alternative = "greater",
  delta = 0,
  sig_level = 0.01)

comparacao_medias

## # A tibble: 1 x 6
## statistic p_value critical_value critical_region delta alternative
### # A tibble: 1 x 6
## statistic p_value critical_value critical_region delta alternative</pre>
```

Ao nível de signifiância 1%, não tem evidência para afirmar que carros automáticos são mais eficientes.

## 1 -3.77 0.999 2.55 (2.548, Inf)

Primeiro precisamos verificar se os desvios padrões são iguais para duas populações.

- Variável aleatória: comprimento de pétala
- População 1: espécie setosa (especies ==
   'setosa')
- População 2: espécie versicolor (especies == 'versicolor')

```
dados <- read xlsx("../dados/brutos/iris.xlsx")</pre>
iris setosa <- dados |> filter(especies == "setosa")
iris versicolor <- dados |> filter(especies == "versicol
comparacao var <- ht 2pop var (
  iris_setosa$comprimento_petala,
  iris_versicolor$comprimento_petala)
comparacao var
## # A tibble 2 x 7
## statistic p_value critical_vale ratio alternative lower_ci upper_ci
##
     <dbl> <dbl>
                     <dbl> <dbl> <chr>
                                        <db1>
                                                <dbl>
## 1 0.137 1.03e-10
                     0.567 1 two.sided
                                        0.137 0.137
## 2 0.137 1.03e-10
                      1.76 1 two.sided
                                        0.137 0.137
```

Ao nível de significância 5%, as variâncias dos comprimentos de pétalas para as duas espécies são diferentes.

Quando sabemos que as variâncias populacionais são diferentes

- População 1:  $N(\mu_1, \sigma)$
- População 1:  $N(\mu_2, \sigma)$
- Teste de Hipóteses envolvendo  $\mu_1$  e  $\mu_2$

No pacote statBasics: ht\_2pop\_mean com argumento var\_equal = F (valor padrão).

#### Testes de hipóteses deste curso:

- Teste bilateral:  $H_1: \mu_1 \mu_2 = \Delta_0$ 
  - alternative = 'two.sided'- valor padrão
- Teste unilateral à esquerda:  $H_1$ :  $\mu_1 \mu_2 < \Delta_0$ 
  - alternative = 'less' (Atenção para ordem das populações)
- Teste unilateral à direita:  $H_1: \mu_1 \mu_2 > \Delta_0$ 
  - alternative = 'greater' (Atenção para ordem das populações)

Especificamos delta fornecendo  $\Delta_0 = \mu_1 - \mu_2$ . Valor padrão: delta = 0 (neste caso, estamos testando a igualdade).

Existe diferença entre os comprimentos médios de pétalas das espécies setosa e versicolor ao nível de significância 5%?

```
comparacao_medias_iris <- ht_2pop_mean(
  iris_setosa$comprimento_petala,
  iris_versicolor$comprimento_petala,
  delta = 0,
  var_equal = T,
  alternative = "two.sided",
  sig_level = 0.05)</pre>
```

Ao nível de significância 5%, os comprimentos médios de pétalas para as espécies setosa e versicolor são diferentes.

- **População 1:** Bernoulli $(p_1)$
- População 1: Bernoulli(p<sub>2</sub>)
- Teste de Hipóteses envolvendo p<sub>1</sub> e p<sub>2</sub>

No pacote statBasics: ht\_2pop\_prop.

#### Duas formas de realizar este Teste de Hipóteses:

- Primeira forma: usando dois vetores de 1 e 0.
- Segunda forma: usando número de sucessos e tamanhos das amostras das duas populações.

#### Testes de hipóteses deste curso:

- Teste bilateral:  $H_1: p_1-p_2=\Delta_0$ 
  - alternative = 'two.sided'- valor padrão
- Teste unilateral à esquerda:  $H_1: p_1 p_2 < \Delta_0$ 
  - alternative = 'less' (Atenção para ordem das populações)
- Teste unilateral à direita:  $H_1: p_1 p_2 > \Delta_0$ 
  - alternative = 'greater' (Atenção para ordem das populações)

Especificamos delta fornecendo  $\Delta_0=p_1-p_2$ . Valor padrão: delta = 0 (neste caso, estamos testando a igualdade).

No cojunto de crédito.xlsx, a proporção de estudantes é igual entre pessoas brancas e negras no contexto de solicitação de crédito ao nível de significância 1%?

```
dados <- read_xlsx("../dados/brutos/credito.xlsx")
dados_branca <- dados |> filter(raca == "Branca")
dados_negra <- dados |> filter(raca == "Negra")
comparacao_prop <- ht_2pop_prop(
  dados_branca$estudante == "Sim",
  dados_negra$estudante == "Sim",
  alternative = "two.sided", sig_level = 0.01)</pre>
```

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

```
comparacao_prop
```

Ao nível de significância 1%, não temos evidência para afirmar que a proporção de estudantes entre pessoas brancas e negras é diferente.

### Teste t pareado

- Uma mesma observação é mensurada antes e depois de um intervenção
- Desejamos checar se a intervenção produziu efeito

Vamos usar a função t.test com o argumento paired = TRUE.

## Teste t pareado

#### Testes de hipóteses deste curso:

- Teste bilateral:  $H_1: \mu_{antes} \neq \mu_{depois}$ 
  - alternative = 'two.sided'- valor padrão
- Teste unilateral à esquerda:  $H_1$ :  $\mu_{antes} < \mu_{depois}$ 
  - alternative = 'less'(Atenção para ordem)
- Teste unilateral à direita:  $H_1: \mu_{antes} > \mu_{depois}$ 
  - alternative ='greater' (Atenção para ordem)

## Teste t pareado

Existe evidência que combinação de dieta e exercício diminuiu a pressão sanguínea ao nível de significância 5%?

```
dados <- read_xlsx(
   "../dados/brutos/pressao_sanguinea.xlsx")
teste_pressao <- t.test(
   dados$antes_exercicio,
   dados$depois_exercicio,
   alternative = "greater",
   paired = T)</pre>
```

#### teste\_pressao

## mean difference

9.3

##

##

Ao nível de significância 5%, o exercício e a dieta produziram efeito na diminuição da pressão sanguínea.

# Teste de associação para variáveis qualitativas

- Variável 1: variável quantitativa
- Variável 2: variável quantitativa
- Queremos checar se Variável 1 e Variável 2 estão associadas

Usamos a função chisq.test do pacote janitor.

#### Queremos testar as seguintes hipóteses:

- *H*<sub>0</sub>: **não** existe associação entre as duas variáveis
- H<sub>1</sub>: existe associação entre as duas variáveis

O tipo de fundação (fundacao\_tipo) e a condição geral da casa (geral\_condicao) estão associadas, ao nível de significância 1%? (do conjunto de dados casas.xlsx)

# Teste de associação para variáveis qualitativas

```
dados <- read_xlsx("../dados/brutos/casas.xlsx")
dados <- dados |>
    mutate(geral_condicao = fct(
        geral_condicao,
        levels = c(
        "muito ruim", "ruim", "abaixo da média",
        "média", "acima da média", "regular",
        "boa", "muito boa", "excelente")))
```

# Teste de associação para variáveis qualitativas

#### Vamos calcular o coeficient V de Cramer.

```
coef_cramer <- CramerV(
  dados$geral_condicao, dados$fundacao_tipo,
  correct = T, conf.level = 0.95)
coef_cramer

## Cramer V lwr.ci upr.ci</pre>
```

```
## 0.2535774 0.2371296 0.2698504
```

```
teste_associacao <- dados |>
  tabyl(geral_condicao, fundacao_tipo) |>
  chisq.test()
```

```
teste_associacao
```

```
##
## Pearson's Chi-squared test
##
## data: tabyl(dados, geral_condicao, fundacao_tipo)
## X-squared = 980.42, df = 40, p-value < 2.2e-16</pre>
```

Ao nível de significância 1%, a condição geral da casa e o tipo de fundação estão associadas.