

---

# Introdução à Estatística no *software* R

## Inferência Estatística

[www.de.ufpb.br](http://www.de.ufpb.br)

<https://www.youtube.com/estatisticalivre>



ESTATÍSTICA APLICADA  
EM SOFTWARE LIVRE

**UFPB**



Departamento de  
**ESTATÍSTICA**



## Usados para indicar a confiabilidade de uma estimativa

- ▶ **Definição:** A partir da amostra procura-se construir um intervalo de variação,  $\hat{\theta}_1 \leq \theta \leq \hat{\theta}_2$ , com certa probabilidade de conter o verdadeiro parâmetro populacional. Isto é, consiste na fixação de dois valores tais que  $(1 - \alpha)$  seja a probabilidade de que o intervalo, por eles determinado, contenha o verdadeiro valor do parâmetro.
  
- ▶ **Tipos de Intervalos de Confiança:**
  1. IC para Média com Variância Populacional Conhecida;
  2. IC para Média com Variância Populacional Desconhecida (“ $n$ ” grande);
  3. IC para Média com Variância Populacional Desconhecida (“ $n$ ” pequena);
  4. IC para Proporção Populacional.

# IC para Média com $\sigma^2$ Conhecida



- ▶ **Definição:** Assumimos que os valores foram amostrados de forma independente e aleatória de uma população com distribuição  $N(\mu; \sigma^2)$ .

$$IC[\mu; 1 - \alpha] = \left[ \bar{X} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}; \bar{X} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right] \quad (1)$$

- ▶ **Bibliotecas Necessárias (R):**

1. `> library(stats)`
2. `> library(TeachingDemos)`
3. `> library(OneTwoSamples)`

- ▶ **Exemplo:** Um pesquisador está estudando a resistência de um dado material sob determinadas condições. Essa variável é normalmente distribuída com desvio padrão de 2 unidades. Utilizando a amostra: 4,9; 7,0; 8,1; 4,5; 5,6; 6,8; 7,2; 5,7; 6,2 unidades, determine o IC para a resistência média com um nível de confiança de 95%.

# IC para Média com $\sigma^2$ Conhecida



## ► Ler o Banco de Dados:

```
> dadosR <- c(4.9, 7.0, 8.1, 4.5, 5.6,  
              6.8, 7.2, 5.7, 6.2)  
  
> dadosR
```

## ► Resolução 1: Função do R (utilizando função de IC)

```
> interval_estimate1(dadosR, sigma = 2, alpha = 0.05)
```

	mean	df	a	b
1	6.222222	9	4.91558	7.528865

## ► Resolução 2: Função do R (utilizando função de TH)

```
> dados2 <- z.test(dadosR, mean(dadosR), stdev= 2,  
                  conf.level = 0.95)  
  
> dados2$conf.int
```

# IC para Média com $\sigma^2$ Desconhecida (“n” grande)



- ▶ **Definição:** Se desconhecermos “ $\sigma^2$ ”, é possível estimá-la por ponto através de “ $S^2$ ”, baseado em uma amostra aleatória.

$$IC[\mu; 1 - \alpha] = \left[ \bar{X} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}; \bar{X} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right] \quad (2)$$

- ▶ **Exemplo:** Foram realizados testes glicêmicos em pacientes após um jejum de 8 horas. Os resultados são apresentados na tabela abaixo. Encontrar um intervalo de confiança de nível 95% para a média  $\mu$ .

80	117	112	91	100	84	104	80	101	95
77	132	118	102	73	103	140	82	92	120
95	78	88	90	102	121	83	88	107	117

# IC para Média com $\sigma^2$ Desconhecida (“ $n$ ” grande)



## ► Ler o Banco de Dados:

```
> dadosG <- c(80, 117, 112, 91, 100, 84, 104, 80, 101,  
              95, 77, 132, 118, 102, 73, 103, 140, 82,  
              92, 120, 95, 78, 88, 90, 102, 121, 83,  
              88, 107, 117)
```

```
> dadosG
```

```
[1] 80 117 112 91 100 84 104 80 101 95 77  
    132 118 102 73 103 140 82 92 120 95 78  
    88 90 102 121 83 88 107 117
```

# IC para Média com $\sigma^2$ Desconhecida (“n” grande)



## ► Resolução 1: Função do R (utilizando função de IC)

```
> interval_estimate1(dadosG, sd(dadosG), alpha = 0.05)
```

	mean	df	a	b
1	99.06667	30	92.91812	105.2152

## ► Resolução 2: Função do R (utilizando função de TH)

```
> dados2 <- z.test(dadosG, mean(dadosG), sd(dadosG),  
  conf.level = 0.95)
```

```
> dados2$conf.int
```

```
[1] 92.91812 105.21522
```

```
attr(,"conf.level")
```

```
[1] 0.95
```

# IC para Média com $\sigma^2$ Desconhecida (“n” pequena)



- ▶ **Definição:** Torna-se necessário uma correção na distribuição padronizada, que consiste em substituir a distribuição Normal Padrão pela distribuição t-Student.

$$IC[\mu; 1 - \alpha] = \left[ \bar{X} - t_{(n-1); \alpha/2} \frac{S}{\sqrt{n}}; \bar{X} + t_{(n-1); \alpha/2} \frac{S}{\sqrt{n}} \right] \quad (3)$$

- ▶ **Exemplo:** Uma amostra de árvores castanheiras, todas com 8 anos de idade, foi observada em uma floresta. Os diâmetros (em polegadas) das árvores foram medidos à uma altura de 3 pés e os resultados foram registrados: 19,4; 21,4; 22,3; 22,1; 20,1; 23,8; 24,6; 19,9; 21,5; 19,1. Construa um IC de 95% para o verdadeiro diâmetro médio das árvores castanheiras dessa idade na floresta.



# IC para Média com $\sigma^2$ Desconhecida (“n” pequena)



## ► Ler o Banco de Dados:

```
> dadosC <- c(19.4, 21.4, 22.3, 22.1, 20.1,  
              23.8, 24.6, 19.9, 21.5, 19.1)  
> dadosC
```

## ► Resolução 1: Função do R (utilizando função de IC)

```
> interval_estimate1(dadosC, sigma = -1, alpha = 0.05)
```

	mean	df	a	b
1	21.42	9	20.10233	22.73767

## ► Resolução 2: Função do R (utilizando função de TH)

```
> dados2 <- t.test(dadosC, conf.level = 0.95)
```

```
> dados2$conf.int
```



- **Definição:** Consideremos a variável aleatória,  $X$ , que representa a presença (ou não) de determinada característica de uma população.

$$IC[p; 1 - \alpha] = \left[ \hat{p} - Z_{\alpha/2} \sqrt{\frac{\hat{p}\hat{q}}{n}}; \hat{p} + Z_{\alpha/2} \sqrt{\frac{\hat{p}\hat{q}}{n}} \right] \quad (4)$$

onde,  $\hat{p} = x/n$  (proporção amostral) e  $\hat{q} = 1 - \hat{p}$ .

- **Exemplo:** Para avaliar a taxa de desemprego em uma cidade, coletou-se uma amostra aleatória de 1000 habitantes em idade de trabalho e observou-se que 87 eram desempregados. Estimar a porcentagem de desempregados em toda cidade através de um intervalo de 95% de confiança.

# IC para a Proporção Populacional



## ► Resolução 1: Função do R (utilizando função de IC)

```
> x <- 87
```

```
> n <- 1000
```

```
# Criar um vetor de 0s e 1s
```

```
> vetor <- c(rep(1,x), rep(0,n-x))
```

```
> vetor
```

```
# Utilizar a função do Intervalo de Confiança
```

```
> interval_estimate1(vetor, sigma = sd(vetor)),  
                    alpha = 0.05)
```

	mean	df	a	b
1	0.087	1000	0.06952326	0.1044767



## ► Resolução 2: Função do R (utilizando função de TH)

```
# Criar um vetor de 0s e 1s (idem Resolução 2)
```

```
# Utilizar a função do Teste de Hipótese
```

```
> teste <- z.test(vetor, stdev = sd(vetor))
```

```
> teste$conf.int
```

```
[1] 0.06952326    0.10447674
```

```
attr(,"conf.level")
```

```
[1] 0.95
```

**Análise:** Conclui-se, ao nível de confiança de 95%, que a proporção de desempregados na cidade, em idade de trabalho, é de no mínimo 6,95% e de no máximo 10,45%.

# TESTES DE HIPÓTESES (TH)



## Processo de Decisão Estatística

- ▶ **Objetivo:** Verificar, através de amostras aleatórias, se hipóteses a respeito de parâmetros populacionais são ou não verdadeiras.
- ▶ **Conceitos Básicos:**
  - 1. Hipótese Estatística
    - ▶ 1.1. Testes Paramétricos
    - ▶ 1.2. Testes Não-Paramétricos
  - 2. Tipos de Hipóteses
    - ▶ 2.1.  $H_0$ : Hipótese Nula
    - ▶ 2.2.  $H_1$ : Hipótese Alternativa
  - 3. Tipos de Erros
    - ▶ 3.1. Erro Tipo I =  $P(\text{Rejeitar } H_0 | H_0 \text{ é Verdadeira})$
    - ▶ 3.2. Erro Tipo II =  $P(\text{Aceitar } H_0 | H_0 \text{ é Falsa})$
  - 4. Tipos de Testes
    - ▶ 4.1. Bilateral
    - ▶ 4.2. Unilateral (à esquerda ou à direita)
  - 5. Estatística de Teste e Região Crítica

# TESTES DE HIPÓTESES PARA UMA POPULAÇÃO



## ► Procedimentos para a Construção dos Testes de Hipóteses:

1. Enunciar as hipóteses:

► **1.1. Bilateral:**  $H_0 : \theta = \theta_0$  versus  $H_1 : \theta \neq \theta_0$

► **1.2. Unilateral à Esquerda:**  $H_0 : \theta = \theta_0$  versus  $H_1 : \theta < \theta_0$

► **1.3. Unilateral à Direita:**  $H_0 : \theta = \theta_0$  versus  $H_1 : \theta > \theta_0$

2. Fixar o nível de significância:  $\alpha$ .

3. Identificar a variável de teste e sua respectiva distribuição de probabilidade: Normal Padrão ( $Z$ ) ou t-Student ( $t$ ).

4. Determinar a região de rejeição.

5. Calcular o valor observado da estatística do teste:

$$Z_{cal} = \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}}$$

ou

$$t_{cal} = \frac{\bar{X} - \mu_0}{S / \sqrt{n}}$$

6. Concluir pela rejeição, ou não, da hipótese  $H_0$ .

# TESTES DE HIPÓTESES



## P-valor

- ▶ Mede a probabilidade de que você tenha tirado seus resultados amostrais de uma  $H_0$  verdadeira.
- ▶ É definido como a probabilidade de se observar um valor da estatística de teste maior ou igual ao encontrado.
- ▶ Quanto mais longe sua estatística de teste estiver com relação às extremidades da distribuição Normal Padrão, menor será o p-valor. Portanto, mais evidências você terá contra a veracidade de  $H_0$ .
- ▶ É uma medida de quanta evidência você tem contra a hipótese nula (quanto **menor** for o p-valor, maior é a evidência para rejeitar  $H_0$ ).
- ▶ Deve-se combinar o p-valor com o nível de significância para tomar a decisão sobre um dado teste de hipótese.

# TESTES DE HIPÓTESES



## P-valor

A tabela seguinte fornece uma interpretação razoável dos P-valores:

P-valor (P)	Interpretação
$P < 0,01$	Evidência muito forte contra $H_0$
$0,01 \leq P < 0,05$	Evidência moderada contra $H_0$
$0,05 \leq P < 0,10$	Evidência sugestiva contra $H_0$
$P \geq 0,10$	Pouca ou nenhuma evidência real contra $H_0$

Tradicionalmente, o valor de corte para rejeitar  $H_0$  é de 0,05, o que significa que, quando não há nenhuma diferença, um valor tão extremo para a estatística de teste é esperado em menos de 5% das vezes.



# Teste de Hipótese para Média com $\sigma^2$ Conhecida



- ▶ **Exemplo:** Seja uma amostra aleatória  $X$  com 1000 observações. Ao nível de significância de 5%, testar se a média populacional é igual a 7 supondo desvio padrão populacional conhecido igual a 2.

- ▶ **Bibliotecas Necessárias aos TH:**

```
> library(TeachingDemos)
```

```
> library(stats)
```

- ▶ **Gerar uma Amostra Aleatória:**

```
> set.seed(08112019)      # Gerar a mesma amostra
```

```
> amostra <- rnorm(1000, mean = 5, sd = 2)
```

# Teste de Hipótese para Média com $\sigma^2$ Conhecida



## ► Resolução: Função do R

```
> z.test(amostra, mu = 7, stdev = 2, conf.level = 0.95)
```

One Sample z-test

```
data:  amostra
```

```
z = -30.353,    n = 1.0000e+03,  Std. Dev. = 2.0000e+00  
Std. Dev. of the sample mean = 6.3246e-02  
p-value < 2.2e-16
```

```
alternative hypothesis: true mean is not equal to 7
```

```
95 percent confidence interval:  
    4.956351    5.204269
```

```
sample estimates:  
mean of amostra  
    5.08031
```

# Teste de Hipótese para Média com $\sigma^2$ Conhecida



## ► Resolução (continuação):

1. Hipóteses:  $H_0 : \mu = 7$  versus  $H_1 : \mu \neq 7$
2. Estatística de Teste:  $Z_{cal} = -30,353$  para  $\alpha = 5\%$
3. Média amostral:  $\bar{X} = 5,0831$ . Desvio Padrão:  $S = 0,0632$
4. Intervalo de Confiança:  $IC[\mu; 95\%] = [4,9564; 5,2043]$
5. P-valor  $< 2,2e - 16$  (menor que 1%)
6. Rejeita-se  $H_0$  ao nível de significância de 5%, ou seja, há evidência suficiente para afirmar que a média populacional é diferente de 7.

# TH para Média com $\sigma^2$ Desconhecida (“n” grande)



- ▶ **Exemplo:** Um pacote de confeitos da marca M&M indica no rótulo um conteúdo de 1498 confeitos com o peso total de 1361g, de modo que o peso médio de cada confeito é 0,9085 g. Em um teste para determinar se o consumidor está sendo prejudicado, seleciona-se uma amostra aleatória de confeitos M&M. Execute o teste ao nível de significância de 5% e verifique a situação do consumidor.

- ▶ **Resolução:**

Importar os registros referentes ao BANCO DE DADOS M&M's

```
> dados <- data.frame(mem)      # Gerar um data frame
```

```
> attach(dados)      # Permite que o R leia as variáveis  
                      contidas no objeto ‘dados’
```

```
> str(dados)
```

# TH para Média com $\sigma^2$ Desconhecida (“n” grande)



## ► Resolução (continuação):

```
# Concatenando os dados
```

```
> mems <- c(Vermelha, Laranja, Amarela, Marrom,  
            Azul, Verde)
```

```
> mems
```

```
[1] 0.751 0.841 0.856 0.799 0.966 0.859 0.857 0.942 ...  
[20] NA NA NA NA NA NA NA NA ...  
[39] 0.793 0.977 0.850 0.830 0.856 0.842 0.778 0.786 ...  
[58] 0.784 0.824 0.858 0.848 0.851 NA NA NA ...  
[77] NA NA NA NA NA 0.696 0.876 0.855 ...  
[96] NA NA NA NA NA NA NA NA ...
```

# TH para Média com $\sigma^2$ Desconhecida (“n” grande)



## ► Resolução (continuação):

```
> length(mems)                # Tamanho do vetor
```

```
[1] 162
```

```
> mems2 <- na.omit(mems)      # Retirar os NA's
```

```
> mems2
```

```
[1] 0.751 0.841 0.856 0.799 0.966 0.859 0.857 0.942 ...  
[20] 0.859 0.838 0.863 0.888 0.925 0.793 0.977 0.850 ...  
[39] 0.883 0.769 0.859 0.784 0.824 0.858 0.848 0.851 ...  
[58] 0.854 0.810 0.858 0.818 0.868 0.803 0.932 0.842 ...  
[77] 0.825 0.869 0.912 0.887 0.886 0.925 0.914 0.881 ...  
[96] 0.778 0.814 0.791 0.810 0.881
```

# TH para Média com $\sigma^2$ Desconhecida (“n” grande)



## ► Resolução (continuação):

```
attr(,"na.action")
[1] 14 15 16 17 18 19 20 21 22 23 24 25 26
    27 53 54 63 64 65 66 67 68 69 70 71 72
[30] 76 77 78 79 80 81 90 91 92 93 94 95 96
    97 98 99 100 101 102 103 104 105 106 107 108 155
[59] 159 160 161 162
```

```
attr(,"class")
[1] "omit"
```

```
> length(mems2)      # Verificar o tamanho da amostra
```

```
[1] 100
```

# TH para Média com $\sigma^2$ Desconhecida (“n” grande)



## ► Resolução (continuação):

```
> testemems2 <- z.test(mems2, mu = 0.9085, sd(mems2),  
                        alternative="less")  
  
> testemems2
```

### One Sample z-test

```
data: mems2  
z = -10.042, n = 1.0000e+02, Std. Dev. = 5.1794e-02,  
Std.Dev. of the sample mean = 5.1794e-03,  
p-value < 2.2e-16  
alternative hypothesis: true mean is less than 0.9085  
95 percent confidence interval:  
      -Inf      0.8650094  
sample estimates:  
mean of mems2  
      0.85649
```



# TH para Média com $\sigma^2$ Desconhecida (“n” grande)



## ► Resolução (continuação):

1. Hipóteses:  $H_0 : \mu = 0,9085 \text{ g}$  versus  $H_1 : \mu < 0,9085 \text{ g}$
2. Estatística de Teste:  $Z_{cal} = -10,042$  para  $\alpha = 5\%$
3. Média amostral:  $\bar{X} = 0,85649$ . Desvio Padrão:  $S = 0,00518$
4. Intervalo de Confiança:  $IC[\mu; 95\%] = [-\infty; 0,8650]$
5. P-valor  $< 2,2e - 16$  (menor que 1%)
6. Rejeita-se  $H_0$  ao nível de significância de 5%, ou seja, há evidência suficiente para desconfiar que o consumidor está sendo prejudicado visto que o peso dos confeitos é menor do que o indicado no rótulo, pelo fabricante.

# TH para Média c/ $\sigma^2$ Desconhecida (“n” pequeno)



- ▶ **Exemplo:** Os valores relacionados a seguir são cargas axiais (em libras) de uma amostra de sete latas de alumínio de 12 oz. A carga axial de uma lata é o peso máximo que seus lados podem suportar, e deve ser superior a 165 libras, porque esta é a pressão máxima aplicada quando se fixa a tampa no lugar. Ao nível de significância de 0,01, teste a afirmação do engenheiro supervisor de que esta amostra provém de uma população com média superior a 165 libras.

270    273    258    204    254    228    282

# TH para Média c/ $\sigma^2$ Desconhecida (“n” pequeno)



## ► Resolução:

```
> latas <- c(270, 273, 258, 204, 254, 228, 282)
```

```
> latas
```

```
[1] 270 273 258 204 254 228 282
```

```
> testelatas <- t.test(latas, mu = 165,  
                        alternative = "greater")
```

# TH para Média c/ $\sigma^2$ Desconhecida (“n” pequeno)



## ► Resolução (continuação):

```
> testelatas
```

```
One Sample t-test
```

```
data:  latas
```

```
t = 8.3984,    df = 6,    p-value = 7.761e-05
```

```
alternative hypothesis: true mean is greater than 165
```

```
95 percent confidence interval:
```

```
232.4193      Inf
```

```
sample estimates:
```

```
mean of x:
```

```
252.7143
```

# TH para Média c/ $\sigma^2$ Desconhecida (“n” pequeno)



## ► Resolução (continuação):

1. Hipóteses:  $H_0 : \mu = 165$  versus  $H_1 : \mu > 165$
2. Estatística de Teste:  $t_{cal} = 8,3984$  para  $\alpha = 1\%$
3. Graus de Liberdade:  $(n - 1) = 6$  graus de liberdade
4. Média amostral:  $\bar{X} = 252,7143$
5. P-valor =  $7,761e-05 = 0,00007761$  (menor que 1%)
6. Rejeita-se  $H_0$  ao nível de significância de 1%, ou seja, há evidência suficiente para apoiar a afirmação do supervisor de que a amostra provém de uma população com média superior às 165 libras desejadas.

# Teste de Hipótese para a Proporção Populacional



- ▶ **Exemplo:** Em um estudo da eficácia do air-bag em automóveis, constatou-se que, em 821 colisões de carros de tamanho médio equipados com air-bag, 46 colisões resultaram em hospitalização do motorista. Ao nível de significância de 0,01, teste a afirmação de que a taxa de hospitalização nos casos de air-bag é inferior à taxa de 7,8% para colisões de carros de tamanho médio equipados com cintos automáticos de segurança.

- ▶ **Resolução 1:**

```
> x <- 46  
> n <- 821  
> pest <- x/n
```

```
# criar vetor de 0s e 1s  
> vetor <- c(rep(1,x), rep(0,n-x))
```

# Teste de Hipótese para a Proporção Populacional

► **Resolução 1 (continuação):**

# Teste de Hipótese para a Proporção Populacional



## ► Resolução 1 (continuação):

1. Hipóteses:  $H_0 : p = 0,078$  versus  $H_1 : p < 0,078$
2. Proporção amostral:  $p_{est} = 0,05602923$
3. Estatística de Teste:  $Z_{cal} = -2,7357$
4. P-valor = 0,003113 (menor que 5%)
5. Rejeita-se  $H_0$  ao nível de significância de 1%, ou seja, há evidência suficiente para apoiar a afirmação de que, para colisões de carros de tamanho médio, a taxa de hospitalização, no caso de haver o air-bag, é inferior à taxa de 7,8% verificada no caso de cintos de segurança automáticos.



# Teste de Hipótese para a Proporção Populacional



**IMPORTANTE:** Uma outra forma de testar a proporção populacional é utilizando a função **prop.test**, em que, sob a hipótese nula, a estatística de teste segue uma distribuição  $\chi^2$ .

## ► Resolução 2:

```
> prop.test(x, n, p = 0.078, alternative = "less")
```

```
1-sample proportions test with continuity correction
```

```
data:  x out of n, null probability 0.078
```

```
X-squared = 5.2094,    df = 1,    p-value = 0.01123
```

```
alternative hypothesis: true p is less than 0.078
```

```
95 percent confidence interval:
```

```
0.00000000    0.07142186
```

```
sample estimates P:
```

```
0.05602923
```

# Exercícios I



1. Seja a seguinte amostra:

40,1	45,0	39,1	43,9	45,8	44,2	37,4	44,7	45,2
41,2	40,7	43,1	44,1	42,6	40,6	41,8	42,9	45,8
43,4	45,5	44,8	42,3	40,4	41,9	42,1	44,4	43,7
43,9	42,6	45,5	41,5	45,2	43,6	42,8	43,3	45,7

Ao nível de  $\alpha = 0,05$  testar:  $H_0 : \mu \leq 42$  contra  $H_1 : \mu > 42$ .

2. O tempo de vida das lâmpadas da marca X tem distribuição aproximadamente normal com a seguinte amostra:

1200	1100	900	1250	1300	1290	1100	1060
1180	1120	1160	1140	1190	1110	1100	1220

Ao nível de  $\alpha = 0,1$  testar:  $H_0 : \mu \geq 1200$  contra  $H_1 : \mu < 1200$ .

## Exercícios II



3. Está sendo realizado um estudo com um novo medicamento. Mais especificamente, deseja-se estimar se o tempo médio de reação após a utilização desse medicamento é menor do que 5 minutos. Para este estudo, o novo medicamento será ministrado em uma amostra de pacientes e o tempo médio de reação de cada um será registrado. Seja a seguinte amostra:

4.0	3.5	6.1	5.8	5.4	4.4	4.9	3.9	5.1	5.3
4.1	4.2	4.8	4.7	3.8	4.8	5.3	5.5	3.6	3.5
4.7	3.3	3.7	6.3	5.7	3.9	4.6	4.7	4.1	4.3

Ao nível de  $\alpha = 0,05$  testar:  $H_0 : \mu \geq 5$  contra  $H_1 : \mu < 5$ .

## Exercícios III



4. Uma pesquisa com 703 trabalhadores selecionados aleatoriamente mostrou que 61% dos respondentes encontraram emprego através de uma rede de amigos. Ache o valor da estatística de teste para a afirmativa de que a maioria (mais de 50%) dos trabalhadores consegue seus empregos através de tais redes.
  
5. Em um estudo dos leitores de código de barras nas lojas, 1234 itens foram verificados, constatando-se que 20 deles foram cobrados em excesso e 1214 itens não o foram. Use o nível de significância 0,05 para testar a afirmativa de que 1% das vendas é cobrado em excesso.

# TESTES DE HIPÓTESES P/ DUAS POPULAÇÕES



## ► Procedimentos para a Construção dos Testes de Hipótese

### 1. Enunciar as hipóteses:

- **1.1.** Bilateral:  $H_0 : \theta_1 = \theta_2$  versus  $H_1 : \theta_1 \neq \theta_2$
- **1.2.** Unilateral à Esquerda:  $H_0 : \theta_1 = \theta_2$  versus  $H_1 : \theta_1 < \theta_2$
- **1.3.** Unilateral à Direita:  $H_0 : \theta_1 = \theta_2$  versus  $H_1 : \theta_1 > \theta_2$

### 2. Fixar o nível de significância: $\alpha$

### 3. Identificar a variável de teste e sua respectiva distribuição de probabilidade

### 4. Determinar a região de rejeição

### 5. Calcular o valor observado da estatística do teste:

$$\boxed{T_{cal} = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}} \quad \text{ou} \quad \boxed{T_{cal} = \frac{\bar{d}}{S_d/\sqrt{n}}} \quad (\text{amostras pareadas})$$

### 6. Concluir pela rejeição, ou não, da hipótese $H_0$ .

# Teste Hipót. para Duas Populações Independentes



- ▶ **Exemplo:** Duas amostras de 10 alunos, de duas turmas distintas, de um mesmo curso apresentam os seguintes totais de pontos em provas de certa disciplina. Ao nível de 5% de significância testar as hipóteses de que as turmas tenham aproveitamento médio diferentes.

Turma 1: 51 47 75 35 72 84 45 11 52 57

Turma 2: 27 75 49 69 73 63 79 37 84 32

- ▶ **Resolução:**

```
> x = c(51, 47, 75, 35, 72, 84, 45, 11, 52, 57)
> y = c(27, 75, 49, 69, 73, 63, 79, 37, 84, 32)
> t.test(x,y)
```

**IMPORTANTE:** Caso haja a suposição de variâncias populacionais iguais, porém desconhecidas, adicione o argumento **var.equal = TRUE**.

## Teste Hipót. para Duas Populações Independentes

► **Resolução (continuação):**

## Welch Two Sample t-test

```
data:  x and y
```

t = -0.62797, df = 17.998, p-value = 0.5379

```
alternative hypothesis: true difference in means
is not equal to 0
```

95 percent confidence interval:

mean of x      mean of y

52.9

# Teste Hipót. para Duas Populações Independentes



## ► Resolução 1 (continuação):

1. Hipóteses:  $H_0 : \mu_1 = \mu_2$  versus  $H_1 : \mu_1 \neq \mu_2$
2. Médias Amostrais:  $\bar{X}_1 = 52,9$  e  $\bar{X}_2 = 58,8$
3. Estatística de Teste:  $t_{cal} = -0,62797$
4. Graus de Liberdade:  $(n - 2) = 17,998$
5. P-valor = 0,5379
6. Não há motivos para rejeitar  $H_0$  ao nível de significância de 5%, ou seja, há evidência suficiente para afirmar que as turmas apresentam aproveitamento médio similar.



# Teste Hipót. para Duas Populações Dependentes



Também conhecido como Teste de Hipóteses para Duas Populações Dependentes (**ou Pareadas**).

- ▶ **Exemplo:** Um consultor que trabalha para o quartel da Polícia Estadual afirma que as armas de serviço dispararão com uma velocidade de boca maior se o cano estiver adequadamente limpo. Obteve-se uma amostra aleatória de armas de 9 mm, e mediu-se a velocidade de boca (em pés por segundo) de um único tiro de cada arma. Cada arma foi profissionalmente limpa e a velocidade de boca de um segundo tiro (com o mesmo tipo de bala) foi medida. Os dados são apresentados na tabela que segue. Verifique se há alguma evidência de que uma arma limpa dispara com velocidade média de boca maior ao nível de significância de 1%.

Antes: 1505 1419 1504 1494 1510 1506

Depois: 1625 1511 1459 1441 1472 1521

# Teste Hipót. para Duas Populações Dependentes



## ► Resolução:

```
> antes <- c(1505, 1419, 1504, 1494, 1510, 1506)
> depois <- c(1625, 1511, 1459, 1441, 1472, 1521)
> t.test(antes, depois, paired = T,
         alternative = "less")
```

```
Paired t-test
data:  antes and depois
t = -0.49656,    df = 5,    p-value = 0.3203
alternative hypothesis: true difference
in means is less than 0
```

```
95 percent confidence interval:
```

```
    -Inf      46.3796
```

```
sample estimates:
```

```
mean of the differences
```

```
-15.16667
```

# Teste Hipót. para Duas Populações Dependentes



## ► Resolução 1 (continuação):

1. Hipóteses:  $H_0 : \mu_D = 0$  versus  $H_1 : \mu_D < 0$
2. Média das Diferenças Amostrais:  $(\bar{X}_1 - \bar{X}_2) = -15,16667$
3. Estatística de Teste:  $t_{cal} = -0,49656$
4. Graus de Liberdade:  $(n - 1) = 5$
5. P-valor = 0,3203
6. Não há motivos para rejeitar  $H_0$  ao nível de significância de 5%, ou seja, não há evidências de que uma arma limpa dispara com velocidade média de boca maior.

# Teste Hipót. para Duas Populações Dependentes



- ▶ **Exemplo:** De 400 moradores sorteados de uma grande cidade industrial, 300 são favoráveis a um projeto governamental. Já na cidade vizinha, de 130 moradores, cuja principal atividade é o turismo, 80 são favoráveis ao projeto governamental. Ao nível de significância de 5%, você diria que as proporções de habitantes favoráveis nas duas cidades são iguais?
- ▶ **Resolução:**

```
> prop.test(x = c(300, 80), n = c(400, 130))
```

```
2-sample test for equality of proportions with continuity correction
```

```
data: c(300, 80) out of c(400, 130)
```

```
X-squared = 8.111, df = 1, p-value = 0.0044
```

```
alternative hypothesis: two.sided
```

```
95 percent confidence interval:  
0.03573921 0.23349156
```

```
sample estimates:  
prop 1 prop 2  
0.7500000 0.6153846
```

# Teste Hipót. para Duas Populações Dependentes



## ► Resolução 1 (continuação):

1. Hipóteses:  $H_0 : p_1 = p_2$  versus  $H_1 : p_1 \neq p_2$
2. Proporções Amostrais:  $\hat{p}_1 = 0,75$  e  $\hat{p}_2 = 0,62$
3. Estatística de Teste:  $\chi^2_{cal} = 8,111$
4. Graus de Liberdade:  $(n - 1) = 1$
5. P-valor = 0,0044
6. Rejeita-se  $H_0$  ao nível de significância de 5%, ou seja, há evidências de que as proporções de habitantes, nas duas cidades, favoráveis ao projeto governamental são diferentes.

# Teste Qui-Quadrado de Independência



- ▶ **Objetivo:** Buscar evidência estatística de que duas variáveis possuem certo grau de associação.
- ▶ **Procedimento para a Construção do Teste:**

1. Enunciar as hipóteses:

$H_0$ : As variáveis são independentes (não estão associadas)

$H_1$ : As variáveis não são independentes (estão associadas)

2. Determinar o nível de significância ( $\alpha$ ) e a variável de teste:

$\chi^2_{[\alpha, (r-1)(c-1)]}$  (distribuição amostral)

3. Determinar a região crítica

4. Calcular a estatística de teste:

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

5. Concluir pela rejeição ou não da hipótese  $H_0$  e interpretar o teste (tomada de decisão).

# Teste Qui-Quadrado de Independência



## ► Critérios para a Validação do Teste:

1. Os dados devem ser selecionados aleatoriamente;
2. Todas as frequências esperadas devem ser maiores ou iguais a 1;
3. Não mais de 20% das frequências esperadas devem ser inferiores a 5;
4. Supondo-se que as variáveis sejam independentes, o valor esperado de cada célula será:

$$E_{r,c} = \frac{(\text{total na linha})(\text{total na coluna})}{\text{total na amostra}}$$

- **Observação:** Se a hipótese de independência ( $H_0$ ) for verdadeira, o valor da estatística de teste será próximo de zero.

# Teste Qui-Quadrado de Independência



- **Exemplo:** A tabela a seguir, apresenta os registros de uma pesquisa realizada com eleitores quanto a classificação por gênero e a sua identificação partidária. Os indivíduos indicaram se eles se identificaram mais fortemente com o Partido Democrata ou o Republicano ou como Independente. Ao nível de significância de 1%, verifique se a preferência partidária está associada ao gênero dos indivíduos.

Gênero	Democrata	Independente	Republicano	Total
Feminino	762 (703,7)	327 (319,6)	468 (533,7)	1557
Masculino	484 (542,3)	239 (246,4)	477 (411,3)	1200
Total	1246	566	945	2757

<sup>1</sup>

<sup>1</sup>Valores entre parênteses são as frequências esperadas estimadas para  $H_0$ .



# Teste Qui-Quadrado de Independência



## ► Resolução:

```
# Gerar a tabela de contingência (TC)
```

```
> dadosTC <- as.table(rbind(c(762, 327, 468),  
                             c(484, 239, 477)))
```

```
> dadosTC
```

	A	B	C
A	762	327	468
B	484	239	477

```
> dimnames(dadosTC) <- list(genero = c("Feminino",  
                                         "Masculino"),  
                             partido = c("Democrata",  
                                         "Independente",  
                                         "Republicano"))
```

# Teste Qui-Quadrado de Independência



## ► Resolução (continuação):

```
> dimnames(dadosTC)
```

```
$genero
```

```
[1] "Feminino" "Masculino"
```

```
$partido
```

```
[1] "Democrata" "Independente" "Republicano"
```

```
> Xsq <- chisq.test(dadosTC)      # Teste Qui-Quadrado
```

```
> Xsq
```

Pearson's Chi-squared test

```
data: dadosTC
```

```
X-squared = 30.07,    df = 2,    p-value = 2.954e-07
```

# Teste Qui-Quadrado de Independência



## ► Resolução (continuação):

```
# Gerar as frequências observadas (Idem objeto "dadosTC")
```

```
> Xsq$observed
```

	partido		
genero	Democrata	Independente	Republicano
Feminino	762	327	468
Masculino	484	239	477

```
# Gerar as frequências esperadas
```

```
> Xsq$expected
```

	partido		
genero	Democrata	Independente	Republicano
Feminino	703.6714	319.6453	533.6834
Masculino	542.3286	246.3547	411.3166

# Teste Qui-Quadrado de Independência



## ► Resolução (continuação):

### 1. Hipóteses:

$H_0$  : A preferência pelo partido independe do gênero

$H_1$  : A preferência pelo partido não independe do gênero

2. Estatística de Teste:  $\chi^2_{(2)} = 30,07$

3. Graus de Liberdade:  $(r - 1)(c - 1) = (2 - 1)(3 - 1) = 2$

4. P-valor = 2,954e-07 (menor que 1%)

5. Rejeita-se  $H_0$  ao nível de significância de 1%, ou seja, há evidência suficiente para afirmar que a preferência pelo partido está associada ao gênero (as variáveis em estudo não são independentes).

# Exercícios I



1. Deseja-se verificar se duas máquinas produzem peças com a mesma homogeneidade quanto à resistência à tensão. Para isso, sorteamos duas amostras de seis peças de cada máquina e obtivemos as seguintes resistências:

Máquina A: 145 127 136 142 141 137

Máquina B: 143 128 132 138 142 132

Ao nível de significância de 0,05 testar se,  $H_0 : \mu_1 = \mu_2$ .

2. Um teste de 200 adultos e 100 adolescentes mostrou que 60 adultos e 50 jovens eram motoristas descuidados. Use estes dados para testar a afirmação de que a porcentagem dos motoristas adolescentes é maior do que a porcentagem dos motoristas adultos ao nível de significância de 5%.

## Exercícios II



3. Deseja-se comparar dois analistas quanto à precisão na análise de uma certa substância que contém carbono. O analista *A* é bastante experiente, e o analista *B* é novo no serviço, sendo, portanto, de experiência desconhecida. Os resultados obtidos foram os seguintes:

A:	-10	16	-8	9	5	-5	5	-11	25	25
B:	-8	-3	20	22	3	5	10	14	-21	8

Ao nível de significância de  $\alpha = 0,10$  teste se há igualdade entre o desempenho médio dos dois analistas, quanto a precisão na análise da substância que contém carbono. Isto é, teste se  $H_0 : \mu_1 = \mu_2$ .



4. Uma empresa que presta serviços de assessoria econômica a outras empresas está interessada em comparar a taxa de reclamações sobre os seus serviços em dois dos seus escritórios em duas cidades diferentes. Suponha que a empresa tenha selecionado aleatoriamente 100 serviços realizados pelo escritório da cidade A e foi constatado que em 12 deles houve algum tipo de reclamação. Já do escritório da cidade B foram selecionados 120 serviços e 18 receberam algum tipo de reclamação. A empresa deseja saber se estes resultados são suficientes para se concluir que os dois escritórios apresentam diferença significativa entre suas taxas de aprovação.

## Exercícios IV



5. Fez-se uma pesquisa para determinar se há restrições, quanto ao gênero, na confiança que o povo deposita na polícia. Os resultados amostrais constam da tabela a seguir. Com o nível de 0,05 de significância, teste a afirmação de que não há tal restrição.

**Tabela:** Número de pessoas entrevistadas quanto a confiança que deposita na polícia segundo o gênero.

Gênero	Confiança na Polícia		
	Muita	Alguma	Muito pouca/Nenhuma
Masculino	115	56	29
Feminino	175	94	31
Total	290	150	60



## Exercícios V



6. Uma pesquisa sobre a qualidade de certo produto foi realizada enviando-se questionários a donas-de-casa através do correio. Aventando-se a possibilidade de que os respondentes voluntários tenham um particular vício de respostas, fizeram-se mais duas tentativas com os não-respondentes. Os resultados estão indicados abaixo. Você acha que existe relação entre a resposta e o número de tentativas?

**Tabela:** Número de donas-de-casa entrevistadas quanto a opinião sobre o produto segundo o número de tentativas.

Opinião sobre o produto	Número de Tentativas		
	1 Tentativa	2 Tentativa	3 Tentativa
Excelente	62	36	12
Satisfatório	84	42	14
Insatisfatório	24	22	24
Total	170	100	50

## Exercícios VI



7. Nicorette é um chiclete que ajuda a fumantes a deixarem de fumar cigarros. A tabela a seguir mostra os resultados de testes feitos para detectar reações negativas. Com nível de 0,05 de significância, teste a afirmação de que o tratamento (remédio ou placebo) é independente da reação (se o paciente experimenta, ou não, irritação na boca ou na garganta). Se o leitor está pensando em recorrer a Nicorette para deixar de fumar, deve se preocupar com aqueles efeitos colaterais?

**Tabela:** Número de pessoas entrevistadas quanto a reação provocada pelo chiclete (Nicorette) segundo o tipo de tratamento.

Reação provocada pelo chiclete	Tratamento	
	Remédio	Placebo
Irritação na boca/garganta	43	35
Nenhuma irritação na boca/garganta	109	118
Total	152	153



8. Um artigo em um jornal discutiu a abertura de um supermercado da rede Whole Food Markets no edifício da Time-Warner na cidade de Nova York. Os dados a seguir compararam os preços de alguns gêneros de primeira necessidade entre os supermercados Whole Food e Fairway, localizado a cerca de 15 quadras do edifício da Time-Warner. Ao nível de significância  $\alpha = 0,01$ , pode-se afirmar que existem evidências de que a média dos preços é mais alta no Whole Foods Market do que no Fairway?

## Exercícios VIII



<i>Item</i>	<i>Whole Foods</i>	<i>Fairway</i>	<i>Item</i>	<i>Whole Foods</i>	<i>Fairway</i>
Leite	2,19	1,35	Trigo	4,99	3,69
Ovos	2,39	1,69	Atum	1,79	1,33
Suco de laranja	2,00	2,49	Maças	1,69	1,49
Salmão	7,99	5,99	Frango	2,19	1,49
Alface	1,98	1,29	Macarrão	1,99	1,59

9. A fim de determinar a eficiência de um medicamento antitêrmico, a temperatura corporal (em graus Celsius) de 20 indivíduos foi medida. Em seguida, foi administrado o medicamento e após uma hora a temperatura foi medida novamente. Os resultados podem ser encontrados na tabela abaixo. Execute o teste de hipóteses ao nível de significância de 1% de modo a avaliar se houve ou não diminuição da temperatura dos indivíduos.

# Exercícios IX



Temperatura			Temperatura		
Indivíduo	Antes	Depois	Indivíduo	Antes	Depois
1	37,5	37,8	11	39,3	38
2	36	36,4	12	37,5	37,1
3	39	37,6	13	38,5	36,6
4	38	37,2	14	39	37,5
5	37,8	36,9	15	36,9	37
6	38,5	37,7	16	37	36,2
7	36,9	36,8	17	38,5	37,6
8	39,4	38,1	18	39	36,8
9	37,2	36,7	19	36,2	36,4
10	38,1	37,3	20	36,8	36,8