



## Testagem de Hipótese Nula

- A lógica da testagem de hipótese nula
- Teste estatístico z de uma média populacional e teste qui-quadrado de uma variância populacional: bilateral e unilateral
- Elementos da decisão estatística: nível de confiança, poder prospectivo, tamanho de efeito populacional e tamanho de amostra
- Valor-p, estatística de teste, tamanho de efeito amostral
- Métodos de planejamento do estudo

$$Z = \frac{\bar{x} - \mu_0}{\frac{\sigma}{\sqrt{n}}}$$

**TESTE Z DE MÉDIA POPULACIONAL**

# Teste z bilateral para uma condição

## Não rejeição de $H_0$

A distribuição da estatura do homem brasileiro de 19 anos de 2016 é normal com desvio-padrão  $\sigma = 7$  cm.

Testar se a média populacional  $\mu$  é igual a  $\mu_0 = 177$  cm hipotetizada.

Quatro participantes desse grupo tiveram suas estaturas medidas, cujos valores em centímetro são: 169, 174, 175, 186.

### Hipóteses

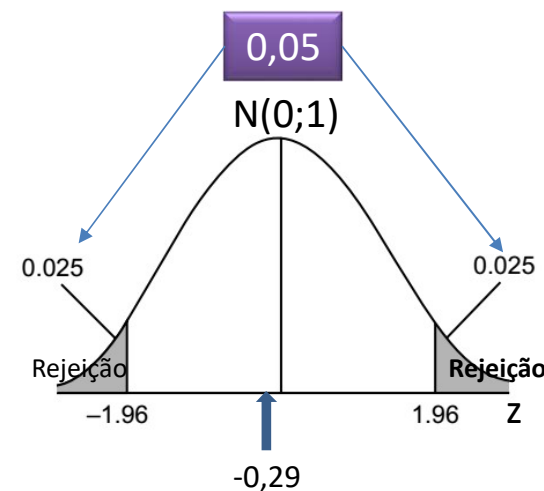
- $H_0: \mu = 177$
- $H_1: \mu \neq 177$  (teste bilateral)

### Estatísticas

- $\bar{X} = (175+186+169+174)/4 = 176$
- $EP = \frac{\sigma}{\sqrt{n}} = 3,5$
- $IC95(\mu) = [176-1,96 \times 3,5; 176+1,96 \times 3,5]$   
 $= [169,14; 182,86]$
- Estatística de teste  $z = \frac{\bar{X}-177}{EP} = \frac{176-177}{3,5} = -0,29$

### Decisão

- *Critério do valor crítico:* Como  $|z| = 0,29 < 1,96$ , não rejeitar  $H_0$  ou
- *Critério do IC95:* Como IC95 contém 177, não rejeitar  $H_0$



### One-Sample Z: ESTATURA

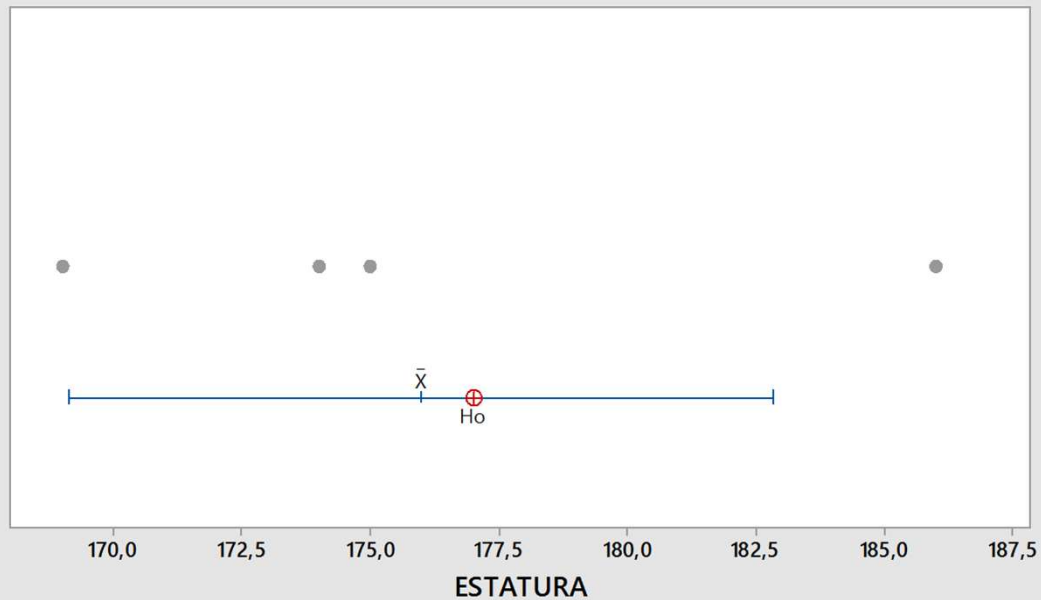
Test of  $\mu = 177$  vs  $\neq 177$

The assumed standard deviation = 7

Variable	N	Mean	StDev	SE Mean	95% CI	Z	P
ESTATURA	4	176,00	7,16	3,50	(169,14; 182,86)	-0,29	0,775

#### Individual Value Plot of ESTATURA

(with Ho and 95% Z-confidence interval for the Mean, and StDev = 7)



MINITAB 17

# Teste z bilateral para uma condição

## *Teste*

- Média populacional  $\mu = 177$  cm hipotetizada

## *Suposições*

- Estatura tem distribuição normal
- Desvio-padrão  $\sigma = 7$  cm conhecido
- $n = 4$  observações independentes: 169, 174, 175, 186
- Teste bilateral
- Nível de confiança adotado de 95% (ou nível de significância de 5%)

## *Hipóteses*

- $H_0: \mu - 177 = 0$  (ausência de efeito)
- $H_1: \mu - 177 \neq 0$

## *Estatísticas*

- $\bar{X} = (169 + 174 + 175 + 186)/4 = 176$
- $EP = \frac{\sigma}{\sqrt{n}} = 3,5$
- $IC95(\mu) = [169,14; 182,86]$
- Estatística de teste  $z = \frac{\bar{X}-177}{EP} = -0,29$

## *Decisão*

- Como  $|z| = 0,29 < 1,96$ , não rejeitar  $H_0$  ou
- Como IC95 contém 177, não rejeitar  $H_0$

# Teste z bilateral para uma condição em R

```
library(BSDA)
estatura <- c(169, 174, 175, 186)
BSDA::z.test(x=estatura, sigma.x=7, mu = 177,
             alternative="two.sided", conf.level=.95)
```

## One-sample z-Test

```
data:  estatura
z = -0.28571, p-value = 0.7751
alternative hypothesis: true mean is not equal to 177
95 percent confidence interval:
 169.1401 182.8599
sample estimates:
mean of x
 176
```

# Teste z bilateral para uma condição

## Rejeição de $H_0$

A distribuição da estatura do homem brasileiro de 19 anos de 2016 é normal com desvio-padrão  $\sigma = 7$  cm.

Testar se a média populacional  $\mu$  é igual a  $\mu_0 = 177$  cm hipotetizada pelo pesquisador.

1.000 participantes desse grupo tiveram suas estaturas medidas, cujos valores em centímetro são: 169, 174, 175, 186, ... .

### Hipóteses

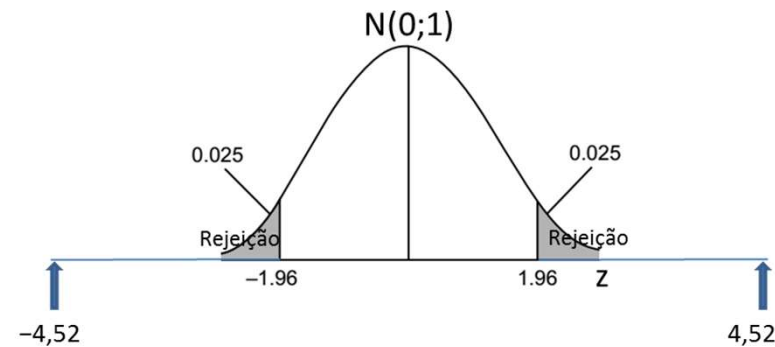
- $H_0: \mu = 177$
- $H_1: \mu \neq 177$

### Estatísticas

- $\bar{X} = (169 + 174 + 175 + 186 + \dots)/1000 = 176$
- $EP = \frac{\sigma}{\sqrt{n}} = 0,22$
- $IC95(\mu) = [175,57; 176,43]$
- Estatística de teste  $z = \frac{\bar{X} - 177}{EP} = -4,52$

### Decisão

- Como  $|z| = 4,52 > 1,96$ , rejeitar  $H_0$  ou
- Como IC95 não contém 177, rejeitar  $H_0$



# Teste z bilateral para uma condição em R

```
library(BSDA)
set.seed(3)
estatura <- rnorm(mean=176, sd=7, n=1000)
BSDA::z.test(x=estatura, sigma.x=7, mu = 177,
             alternative="two.sided", conf.level=.95)
```

## One-sample z-Test

```
data:  estatura
z = -4.3153, p-value = 1.594e-05
alternative hypothesis: true mean is not equal to 177
95 percent confidence interval:
 175.6109 176.4786
sample estimates:
mean of x
 176.0448
```



# Valor-p

- O valor-p é a probabilidade de que a estatística de teste seja igual ou mais extrema que o valor observado na direção prevista pela hipótese alternativa ( $H_1$ ), presumindo que a hipótese nula ( $H_0$ ) é verdadeira.
- AGRESTI, A. & FINLAY, B. (2012) *Métodos estatísticos para as Ciências Sociais*. Porto Alegre: PENSO, p. 171.

# Valor-p

$$\text{p-value} = \frac{\Gamma\left(\frac{(n+1)}{2}\right)}{\sqrt{n \cdot \pi} \Gamma\left(\frac{n}{2}\right)} \int_{-\infty}^t \left(1 + \frac{x^2}{n}\right)^{-\left(\frac{(n+1)}{2}\right)} dx$$

# Teste z bilateral para uma condição

A distribuição da estatura do homem brasileiro de 19 anos de 2016 é normal com desvio-padrão  $\sigma = 7$  cm.

Testar se a média populacional  $\mu$  é igual a  $\mu_0 = 177$  cm hipotetizada.

Quatro participantes desse grupo tiveram suas estaturas medidas, cujos valores em centímetro são: 169, 174, 175, 186.

## Hipóteses

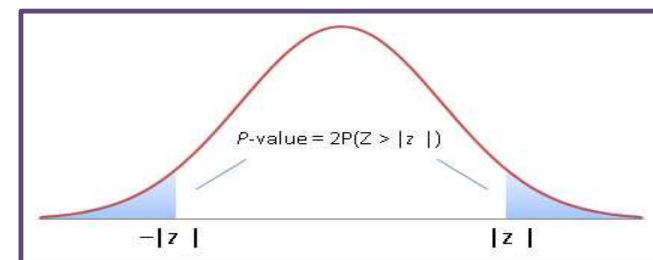
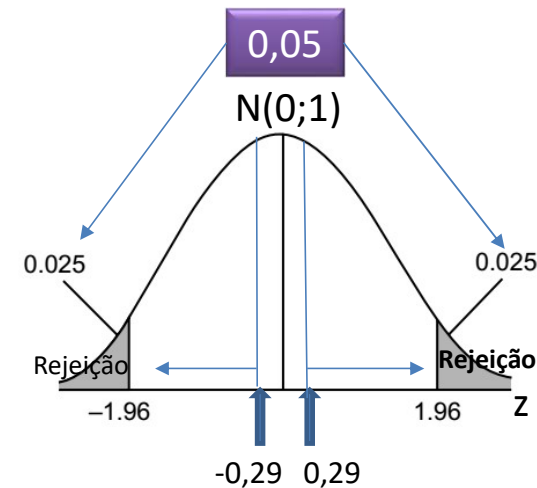
- $H_0: \mu = 177$
- $H_1: \mu \neq 177$

## Estatísticas

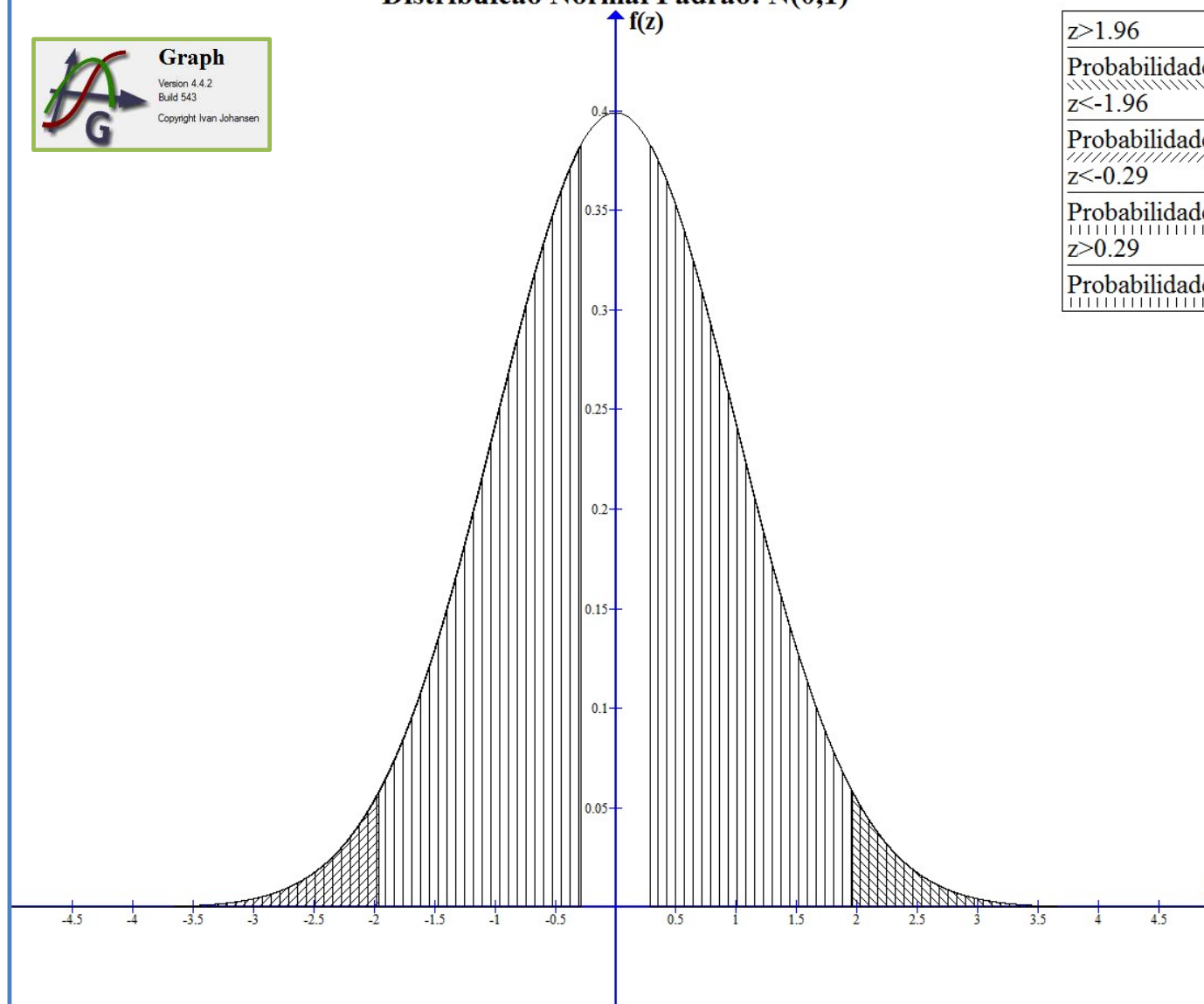
- $\bar{X} = (169 + 174 + 175 + 186)/4 = 176$
- $EP = \frac{\sigma}{\sqrt{n}} = 3,5$
- $IC95(\mu) = [176 - 1,96 \times 3,5; 176 + 1,96 \times 3,5]$   
 $= [169,14; 182,86]$
- Estatística de teste  $z = \frac{\bar{X} - \mu_0}{EP} = \frac{176 - 177}{3,5} = -0,29$

## Decisão

- Como  $|z| = 0,29 < 1,96$ , não rejeitar  $H_0$  ou
- Como IC95 contém 177, não rejeitar  $H_0$  ou
- *Critério do valor-p*: Como a probabilidade de escores-z serem mais extremos que 0,29 e -0,29, i.e., o valor-p bilateral  $= 0,77 = 2 * pnorm(-abs(-0.29))$  é maior que 5%, não rejeitar  $H_0$



# Distribuicao Normal Padrao: N(0,1)



# Teste z bilateral para uma condição em R

```
library(BSDA)
estatura <- c(169, 174, 175, 186)
BSDA::z.test(x=estatura, sigma.x=7, mu = 177,
             alternative="two.sided", conf.level=.95)
```

## One-sample z-Test

```
data:  estatura
z = -0.28571, p-value = 0.7751
alternative hypothesis: true mean is not equal to 177
95 percent confidence interval:
 169.1401 182.8599
sample estimates:
mean of x
 176
```

# Teste z bilateral para uma condição

A distribuição da estatura do homem brasileiro de 19 anos de 2016 é normal com desvio-padrão  $\sigma = 7$  cm.

Testar se a média populacional  $\mu$  é igual a  $\mu_0 = 177$  cm hipotetizada.

Mil participantes desse grupo tiveram suas estaturas medidas, cujos valores em centímetro são: 169, 174, 175, 186, ....

## Hipóteses

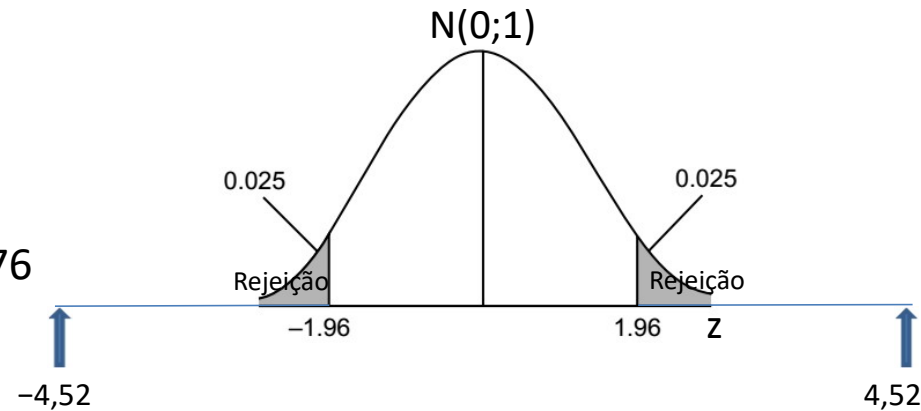
- $H_0: \mu = 177$
- $H_1: \mu \neq 177$

## Estatísticas

- $\bar{X} = (169 + 174 + 175 + 186 + \dots)/1000 = 176$
- $EP = \frac{\sigma}{\sqrt{n}} = 0,22$
- $IC95(\mu) = [175,57; 176,43]$
- Estatística de teste  $z = \frac{\bar{X} - 177}{EP} = -4,52$

## Decisão

- Como  $|z| = 4,52 > 1,96$ , rejeitar  $H_0$  ou
- Como IC95 não contém 177, rejeitar  $H_0$
- Como a probabilidade de escores-z serem mais extremos que -4,52 e 4,52, i.e., o valor-p bilateral =  $6,18E-06 = 2 * pnorm(-abs(-4.52))$  é menor que 5%, rejeitar  $H_0$



# Teste z bilateral para uma condição

## *Teste*

- Média populacional  $\mu = 177$  cm hipotetizada

## *Suposições*

- Estatura tem distribuição normal
- Desvio-padrão  $\sigma = 7$  cm conhecido
- $n=1000$  observações independentes
- Teste bilateral
- Nível de confiança de 95% (ou nível de significância adotado de 5%)

## *Hipóteses*

- $H_0: \mu - 177 = 0$  (ausência de efeito)
- $H_1: \mu - 177 \neq 0$

## *Estatísticas*

- $\bar{X} = (169 + 174 + 175 + 186 + \dots)/1000 = 176$
- $EP = \frac{\sigma}{\sqrt{n}} = 0,22$
- $IC95(\mu) = [175,57; 176,43]$
- Estatística de teste  $z = \frac{\bar{X} - 177}{EP} = -4,52$

## *Decisão*

- Critério do valor crítico da estatística de teste: Como  $|z| = 4,52 > 1,96$ , rejeitar  $H_0$  ou
- Critério do IC95: Como IC95 não contém 177, rejeitar  $H_0$
- Critério do valor-p: Como a probabilidade de escores-z serem mais extremos que -4,52 e 4,52, i.e., o valor-p bilateral =  $6,18E-06 = 2 * pnorm(-abs(-4.52))$  é menor que 5%, rejeitar  $H_0$

# Teste z bilateral para uma condição em R

```
library(BSDA)
set.seed(3)
estatura <- rnorm(mean=176, sd=7, n=1000)
BSDA::z.test(x=estatura, sigma.x=7, mu = 177,
             alternative="two.sided", conf.level=.95)
```

## One-sample z-Test

```
data:  estatura
z = -4.3153, p-value = 1.594e-05
alternative hypothesis: true mean is not equal to 177
95 percent confidence interval:
 175.6109 176.4786
sample estimates:
mean of x
 176.0448
```



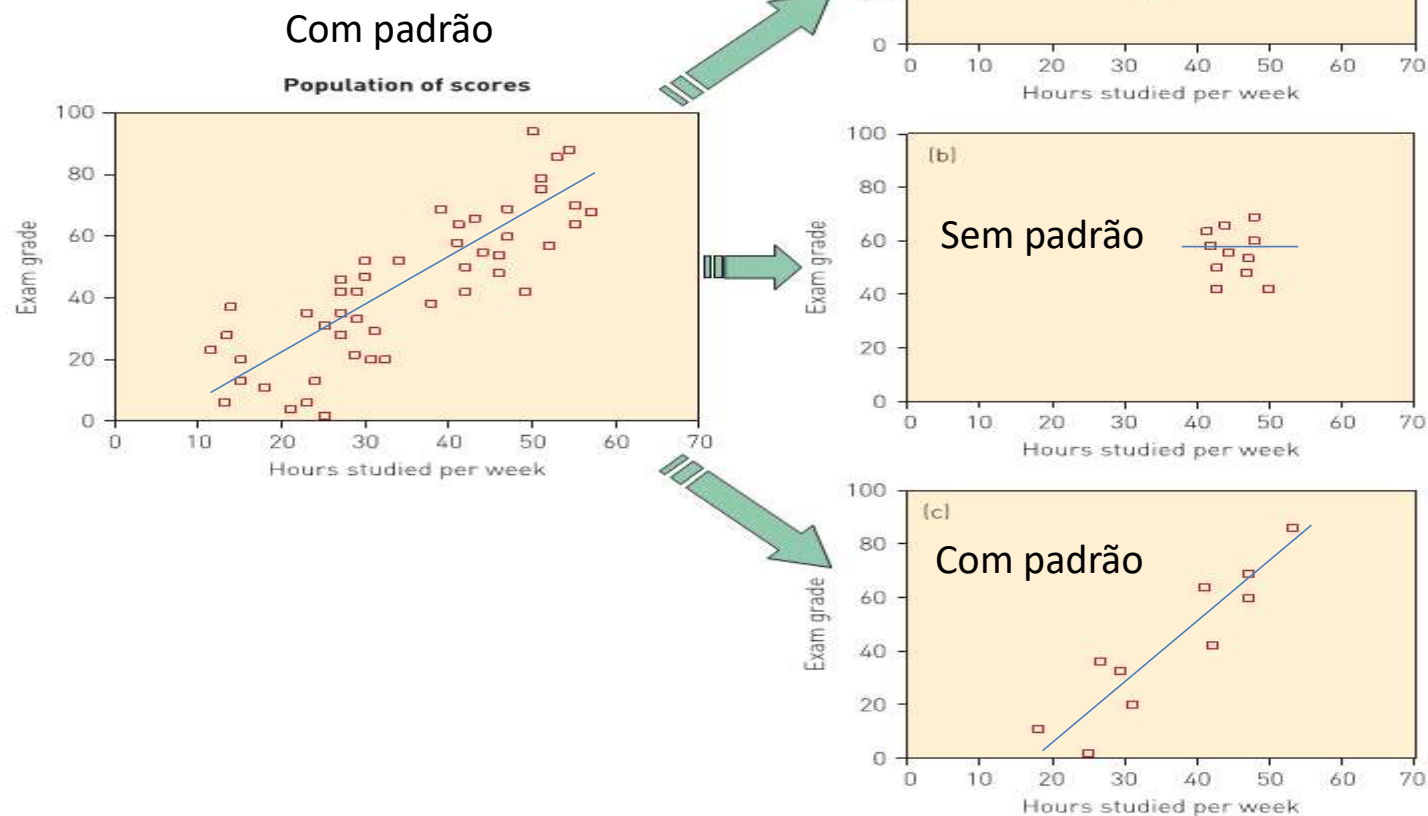
## Erro (Flutuação) Amostral

- Devido ao erro amostral, as amostras que utilizamos podem não refletir de forma fiel a população de onde foram retiradas.

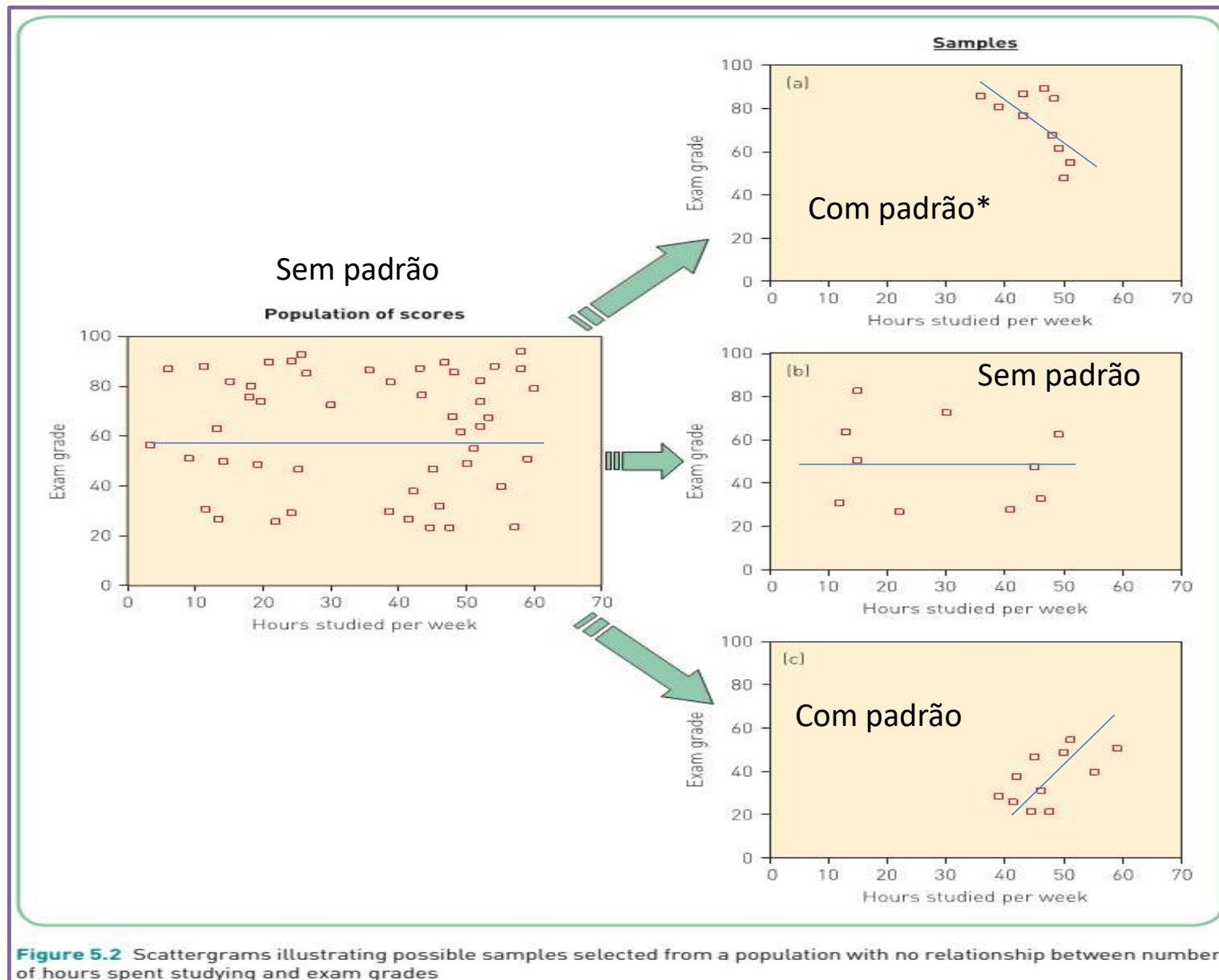
# Conceito do valor-p

- Constitui um dos problemas enfrentados quando conduzimos uma pesquisa o fato de não sabermos qual é o padrão existente na população de interesse.
- De fato, o motivo de realizarmos a pesquisa é, em primeiro lugar, determinar esse padrão.
- Você precisa estar ciente de que, algumas vezes, devido ao erro amostral, obteremos padrões nas amostras que não refletem de forma acurada a população de onde as amostras foram retiradas.
- Assim, precisamos de um algum meio para avaliar a probabilidade de que a amostra selecionada seja um retrato fiel da população.
- Os testes estatísticos nos auxiliam nesta decisão, mas isso ocorre de uma forma não de todo intuitiva.

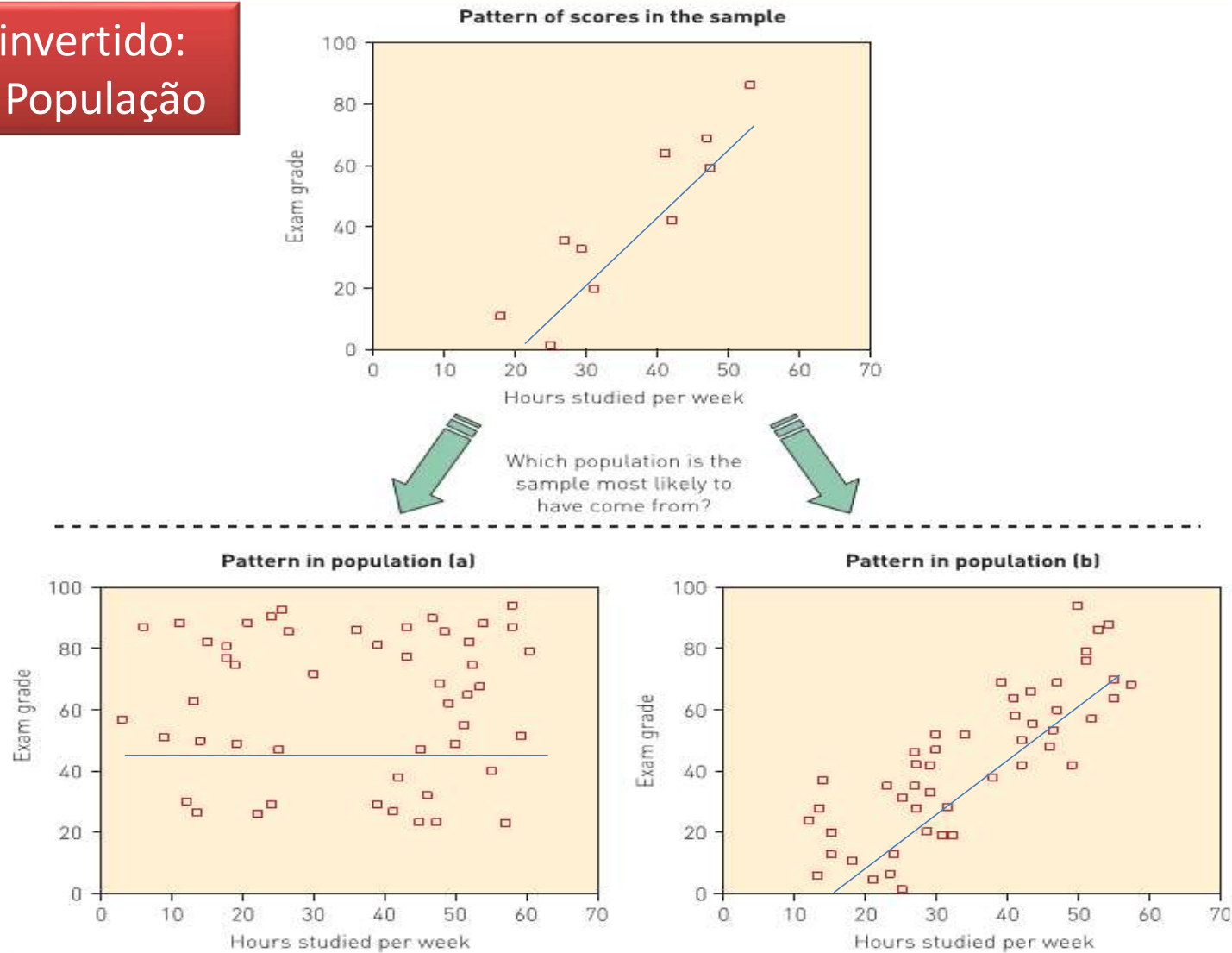
Problema:  
População -> Amostra



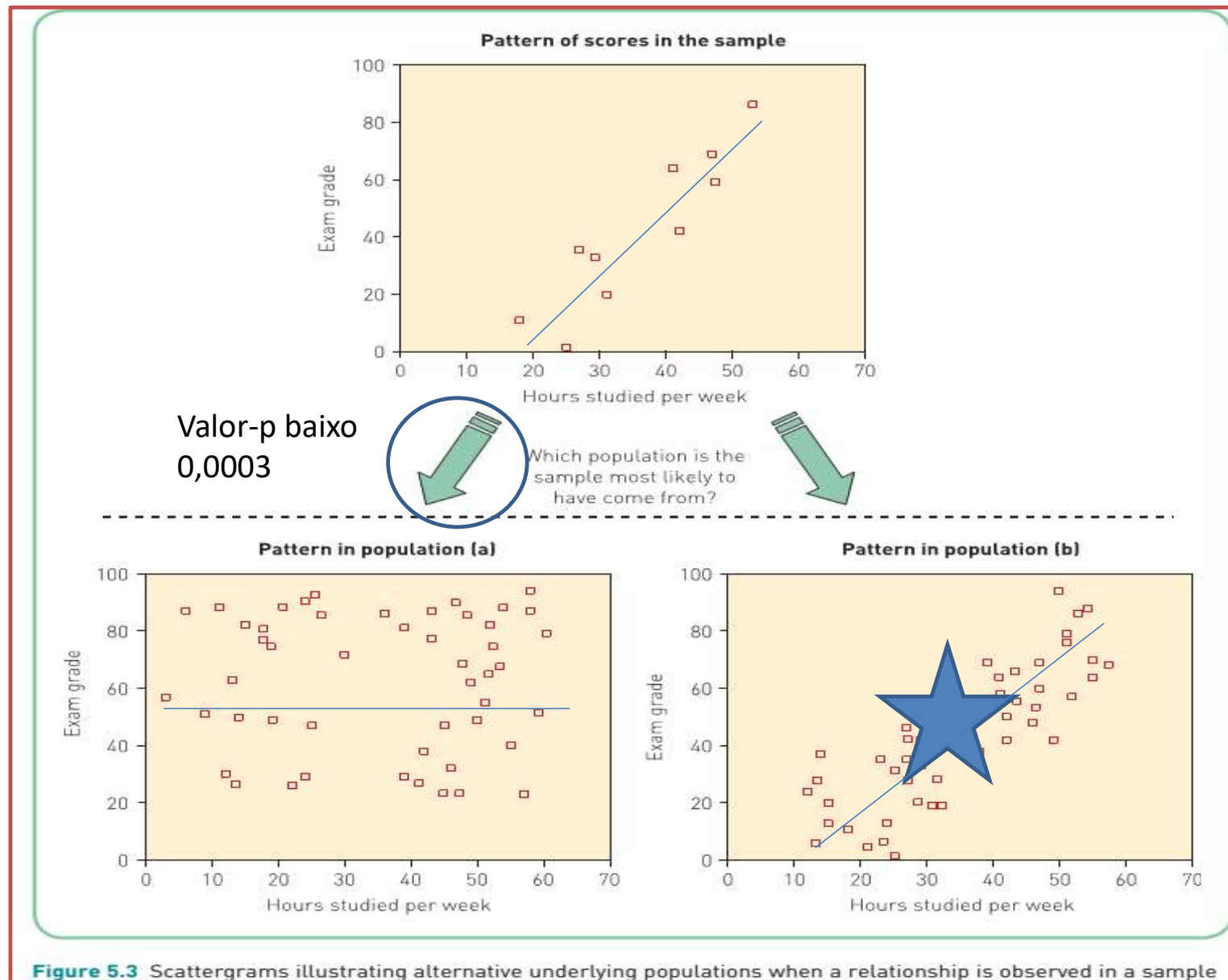
**Figure 5.1** Scattergrams illustrating possible samples selected from a population with a positive relationship between number of hours spent studying and exam grades

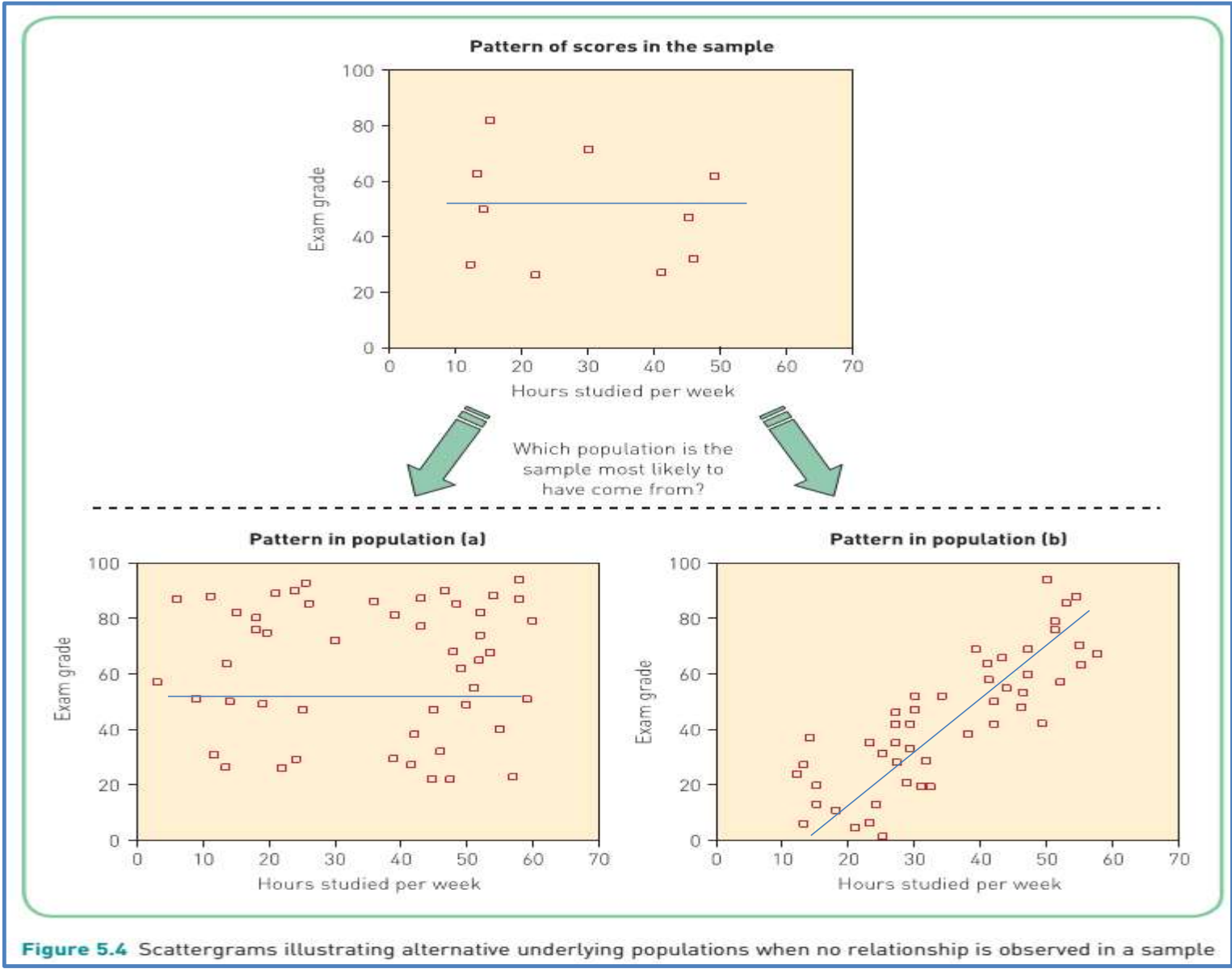


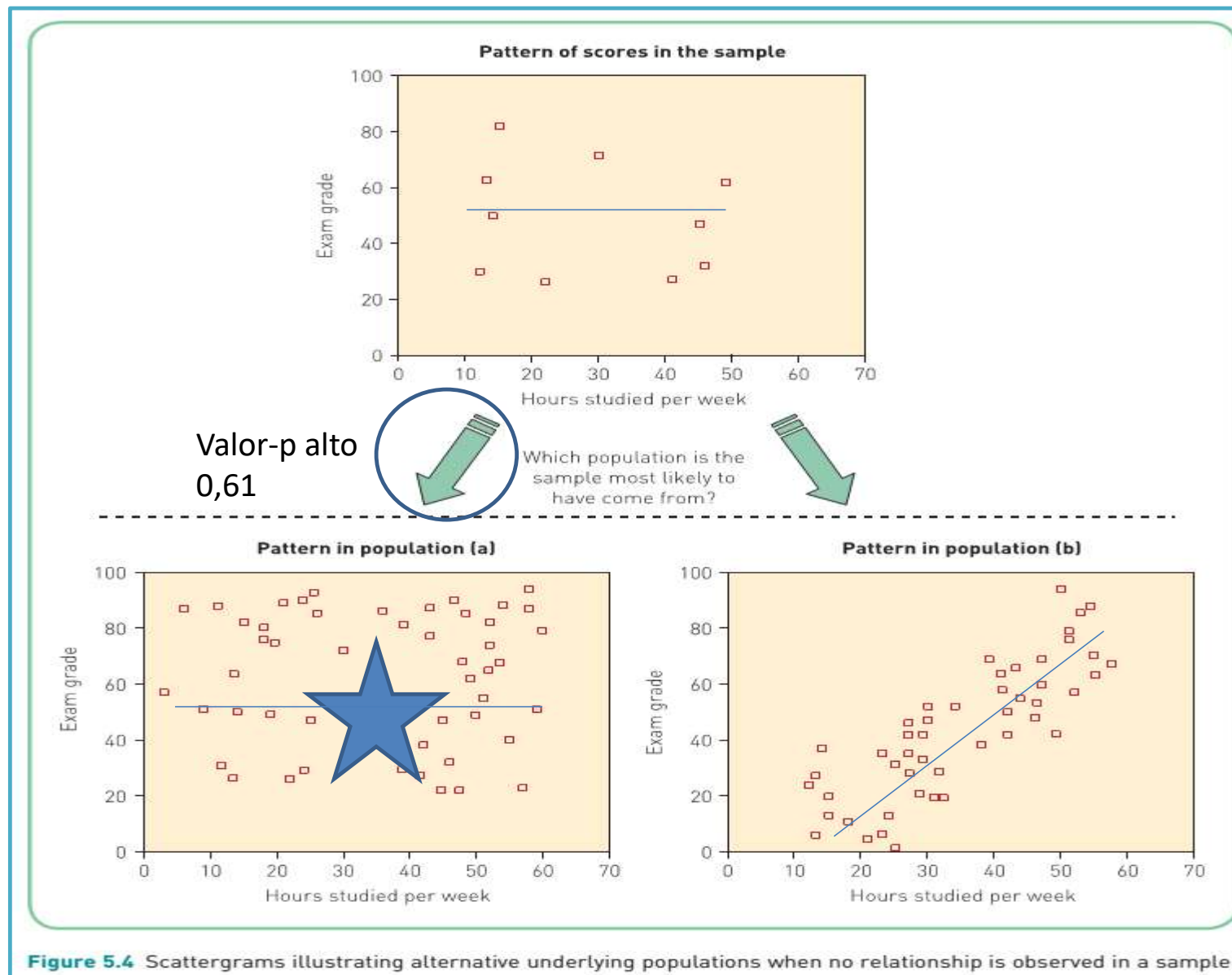
Problema invertido:  
Amostra -> População



**Figure 5.3** Scattergrams illustrating alternative underlying populations when a relationship is observed in a sample









# Efeito

- Correlação entre variáveis
- Diferença entre condições

# Hipótese nula

- Definição de *hipótese nula* ou  $H_0$ 
  - A hipótese nula sempre declara que não existe efeito na população.
- Definição de *hipótese de pesquisa* ou alternativa ou  $H_1$  ou  $H_a$ 
  - A hipótese de pesquisa é a nossa previsão de como grupos específicos podem estar relacionados entre si.
  - De forma alternativa, pode ser nossa previsão de como grupos específicos de participantes podem ser diferentes entre si ou como um grupo de participantes pode ser diferente quando tem um desempenho sob duas ou mais condições experimentais.

# Decisão Estatística VS. Estado da Natureza



## Possible Outcomes of the Decision-Making Process

Researcher's Decision	True State of the World	
	$H_0$ True	$H_0$ False
Reject $H_0$	Type I error $p = \alpha$ = significance level	Correct Decision $p = 1 - \beta$ = Power
Fail to Reject $H_0$	Correct Decision $p = 1 - \alpha$ = confidence level	Type II error $p = \beta$

## Teste de Hipótese Nula & Intervalo de Confiança

- Rejeitar a hipótese nula ao nível de significância adotado,  $\alpha$ , se o valor do parâmetro conjecturado na hipótese nula não pertencer ao intervalo de confiança de  $1 - \alpha$ .
- Os critérios de valor-p e do IC são equivalentes.

# Críticas contra os testes de hipótese nula



- A testagem da hipótese nula é a abordagem dominante na Psicologia e Medicina
- Apesar das críticas à testagem da hipótese nula, isso não significa que tal abordagem deve ser abandonada completamente
- Ao invés disso, devemos ter um entendimento completo de seu significado para podermos nos beneficiar desta tecnologia da decisão
- Além do valor-p, é importante usar o intervalo de confiança e de tamanho de efeito

## Recurring controversies about $P$ values and confidence intervals revisited

ARIS SPANOS<sup>1</sup>

*Department of Economics, Virginia Tech, Blacksburg, Virginia 24061 USA*

### *P value and the large $n$ problem*

A crucial weakness of both the  $P$  value and the N-P error probabilities is the so-called large  $n$  problem: there is always a large enough sample size  $n$  for which any simple null hypothesis.  $H_0: \mu = \mu_0$  will be rejected by a frequentist  $\alpha$ -significance level test; see Lindley (1957).

The large  $n$  constitutes an example of a broader problem known as the *fallacy of rejection*: (mis)interpreting reject  $H_0$  (evidence against  $H_0$ ) as evidence for a particular  $H_1$ ; this can arise when a test has very high power, e.g., large  $n$ . A number of attempts have been made to alleviate the large  $n$  problem, including rules of thumb for decreasing  $\alpha$  as  $n$  increases; see Lehmann (1986). Due to the trade-off between the Type I and II error probabilities, however, any attempt to ameliorate the problem renders the inference susceptible to the reverse fallacy known as the *fallacy of acceptance*: (mis)interpreting accept  $H_0$  (no evidence against  $H_0$ ) as evidence for  $H_0$ ; this can easily arise when a test has very low power; e.g.,  $\alpha$  is tiny or  $n$  is too small.

These fallacies are routinely committed by practitioners in many applied fields. After numerous unsuccessful attempts, Mayo (1996) provided a reasoned answers to these fallacies in the form of a post-data severity assessment.



# Significância prática

- Mesmo efeitos muito pequenos poderão apresentar significância estatística quando o tamanho da amostra for bem grande
- Para determinar a significância prática a melhor abordagem consiste em obter uma medida do tamanho do efeito, sendo que essa medida não depende do tamanho da amostra
  - E.g.: a correlação de Pearson amostral mede a intensidade da associação linear entre duas variáveis quantitativas e não depende do tamanho da amostra



# Interpretação errônea do valor-p

- Muitos pesquisadores sem experiência em estatística (e mesmo aqueles com alguma) equiparam o valor-p com o verdadeira tamanho do efeito, i.e., quanto menor o valor-p, mais forte seria, por exemplo, o relacionamento entre duas variáveis; talvez, de fato, quanto mais forte o relacionamento, mais baixo o valor-p, mas não significa que isso necessariamente ocorrerá
- **O valor-p não é a probabilidade de que a hipótese nula seja verdadeira**; de fato, não sabemos qual é a probabilidade de que a hipótese nula seja verdadeira
- $1 - p$  não é a probabilidade de que a hipótese alternativa seja verdadeira; de fato, não sabemos qual é a probabilidade de que a hipótese alternativa seja verdadeira

# Understanding the Role of *P* Values and Hypothesis Tests in Clinical Research

JAMA Cardiol. doi:10.1001/jamacardio.2016.3312

Published online October 12, 2016.

Daniel B. Mark, MD, MPH; Kerry L. Lee, PhD; Frank E. Harrell Jr, PhD

*P* values and hypothesis testing methods are frequently misused in clinical research. Much of this misuse appears to be owing to the widespread, mistaken belief that they provide simple, reliable, and objective triage tools for separating the true and important from the untrue or unimportant. The primary focus in interpreting therapeutic clinical research data should be on the treatment ("oomph") effect, a metaphorical force that moves patients given an effective treatment to a different clinical state relative to their control counterparts. This effect is assessed using 2 complementary types of statistical measures calculated from the data, namely, effect magnitude or size and precision of the effect size. In a randomized trial, effect size is often summarized using constructs, such as odds ratios, hazard ratios, relative risks, or adverse event rate differences. How large a treatment effect has to be to be consequential is a matter for clinical judgment. The precision of the effect size (conceptually related to the amount of spread in the data) is usually addressed with confidence intervals. *P* values (significance tests) were first proposed as an informal heuristic to help assess how "unexpected" the observed effect size was if the true state of nature was no effect or no difference. Hypothesis testing was a modification of the significance test approach that envisioned controlling the false-positive rate of study results over many (hypothetical) repetitions of the experiment of interest. Both can be helpful but, by themselves, provide only a tunnel vision perspective on study results that ignores the clinical effects the study was conducted to measure.

Table 2. Common Misconceptions About *P* Value

Misconception	Comment
<i>P</i> value equals the probability that the null hypothesis is true.	<i>P</i> value is computed by assuming the null hypothesis is true.
<i>P</i> value equals the probability that the observed effect is due to "the play of chance."	<i>P</i> value is defined as the probability of a difference (effect) as large as that observed or larger if the null hypothesis is true. Even if the difference observed is consistent with a simple chance mechanism, other more complex explanations are also possible, and nothing in <i>P</i> value calculation allows one to conclude that this is the best or most likely explanation for the observed differences.
<i>P</i> value $\leq .05$ means the null hypothesis is false. <i>P</i> value $> .05$ means the null hypothesis is true.	<i>P</i> value is computed assuming the null hypothesis is true. It is not the probability that the null hypothesis is either true or false.
<i>P</i> value $\leq .05$ identifies a clinically or scientifically important difference (effect). <i>P</i> value $> .05$ rules out a clinically or scientifically important difference (effect).	Clinical or scientific importance of study results is a judgment integrating multiple elements, including effect size (expected and observed), precision of estimate of effect size, and knowledge of prior relevant research. At best, <i>P</i> value has a minor role in shaping this judgment.
A small <i>P</i> value indicates study results are reliable and likely to replicate.	<i>P</i> value provides no information about whether a given study result can be reproduced in a second, replication experiment. There are many other factors that must be considered in judging the reliability of study results. Understanding what works in medicine is a process and not the product of any single experiment.



# A Dirty Dozen: Twelve *P*-Value Misconceptions

Steven Goodman

The *P* value is a measure of statistical evidence that appears in virtually all medical research papers. Its interpretation is made extraordinarily difficult because it is not part of any formal system of statistical inference. As a result, the *P* value's inferential meaning is widely and often wildly misconstrued, a fact that has been pointed out in innumerable papers and books appearing since at least the 1940s. This commentary reviews a dozen of these common misinterpretations and explains why each is wrong. It also reviews the possible consequences of these improper understandings or representations of its meaning. Finally, it contrasts the *P* value with its Bayesian counterpart, the Bayes' factor, which has virtually all of the desirable properties of an evidential measure that the *P* value lacks, most notably interpretability. The most serious consequence of this array of *P*-value misconceptions is the false belief that the probability of a conclusion being in error can be calculated from the data in a single experiment without reference to external evidence or the plausibility of the underlying mechanism.

Semin Hematol 45:135-140 © 2008 Elsevier Inc. All rights reserved.

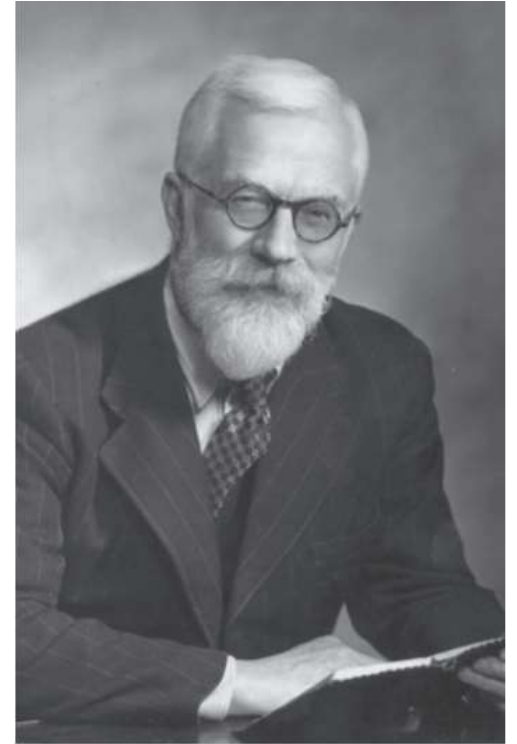
**Table 1** Twelve *P*-Value Misconceptions

1	<i>If <math>P = .05</math>, the null hypothesis has only a 5% chance of being true.</i>
2	<i>A nonsignificant difference (eg, <math>P \geq .05</math>) means there is no difference between groups.</i>
3	<i>A statistically significant finding is clinically important.</i>
4	<i>Studies with <i>P</i> values on opposite sides of .05 are conflicting.</i>
5	<i>Studies with the same <i>P</i> value provide the same evidence against the null hypothesis.</i>
6	<i><math>P = .05</math> means that we have observed data that would occur only 5% of the time under the null hypothesis.</i>
7	<i><math>P = .05</math> and <math>P \leq .05</math> mean the same thing.</i>
8	<i><i>P</i> values are properly written as inequalities (eg, "<math>P \leq .02</math>" when <math>P = .015</math>)</i>
9	<i><math>P = .05</math> means that if you reject the null hypothesis, the probability of a type I error is only 5%.</i>
10	<i>With a <math>P = .05</math> threshold for significance, the chance of a type I error will be 5%.</i>
11	<i>You should use a one-sided <i>P</i> value when you don't care about a result in one direction, or a difference in that direction is impossible.</i>
12	<i>A scientific conclusion or treatment policy should be based on whether or not the <i>P</i> value is significant.</i>

# Replicação

- A replicação é uma das pedras angulares da ciência
- Se você observa um fenômeno uma vez, então pode ter sido por acaso; se o observa duas, três ou mais vezes, pode estar começando a aprender algo sobre o fenômeno estudado
- Se o seu estudo foi o primeiro neste assunto, é sensato que você trate os resultados com certo grau de cautela

# Por que estabelecer $\alpha = 5\%$ ?



The great R. A. Fisher wrote in 1926: "Personally, the writer prefers to set a low standard of significance at the 5 percent point, and ignore entirely all results which fail to reach that level. A scientific fact should be regarded as experimentally established only if a properly designed experiment rarely fails to give this level of significance." (Quoted in Moore, 1979 edition).

It is the fate of a guru that what he sees as a convenient but arbitrary option is taken by followers as written in stone. But it is a philosophy that must be abandoned.

Moore, D. S. (1997) Statistics: Concepts and Controversies, 4<sup>th</sup> edition. New York: Freeman

# Testes unilaterais e bilaterais

- Quando a direção do relacionamento ou da diferença (efeito) é especificada, então o teste é unilateral/unicaudal; caso contrário, é bilateral/bicaudal.
- Em geral (mas nem sempre), se você tiver obtido um valor-p para um teste bilateral e quiser saber o valor mínimo correspondente para o teste unilateral, então:  $p_{\text{uni}} \geq p_{\text{bi}}/2$
- Observe que o que deve ser dobrado ou dividido por 2 não é a estatística de teste (e.g.: z).



## Null Hypothesis ( $H_0$ ) and Alternative Hypothesis ( $H_1$ )

In the above three example research problems, a *treatment* is applied to the members of a sample. In the first example, the treatment is a daily dose of alcohol given to pregnant rats. In the second, rats living in a cool environment is the treatment. In the third example, each person is given a drug for stress. These treatments are being tested to see if they affect the sample members. In the first example research problem, can we conclude that alcohol usage reduces the birth weight of the rats? In the second example, can we conclude that living in a cool environment will increase food consumption? In the third example, does the drug affect response time? In each of the examples the *null hypothesis* is that the treatment has a null or zero effect. We indicate this symbolically in the three examples as:  $H_0: \mu = 20$  grams or the alcohol has no affect on birth weight;  $H_0: \mu = 12$  grams or the cool temperature has no effect on the food consumption; and  $H_0: \mu = 10$  seconds or the drug has no effect on response time. The *alternative hypotheses* are represented as  $H_1: \mu < 20$  grams or the alcohol does reduce the average birth weight;  $H_1: \mu > 12$  grams or the rats consume more food in the cooler environment; and  $H_1: \mu \neq 10$  seconds or the drug affects response time. In each of these situations,  $\mu$  represents the mean after the treatment has been applied.

## Reaching a Decision

What is the probability of getting a sample of size  $n = 50$  with a sample mean of  $M = 18$  if the sample comes from a population with mean  $\mu$  equal to 20? It is now that we turn to the results in Chapter 7 on the sampling distribution of the mean. We assume the null hypothesis to be true and calculate the probability of getting a sample mean of 18 or smaller from a population with  $\mu = 20$ . The distribution of  $M$  is normal with mean = 20 and standard error equal to  $\sigma/\sqrt{n}$  or  $4/\sqrt{50}$ , which equals 0.57 as shown in Figure 8.1.

$H_0: \mu \geq \mu_0$  é **equivalente** a  $H_0: \mu = \mu_0$   
vs.  
 $H_1: \mu < \mu_0$

Demonstração em  
GATÁS, RR (1978) *Elementos de Probabilidade e Inferência*. SP: Atlas, p. 220-3.



```
pnorm(q=18, mean=20, sd=0.57, lower.tail=TRUE)
[1] 0.0002250904
```

TCL (n > 30)

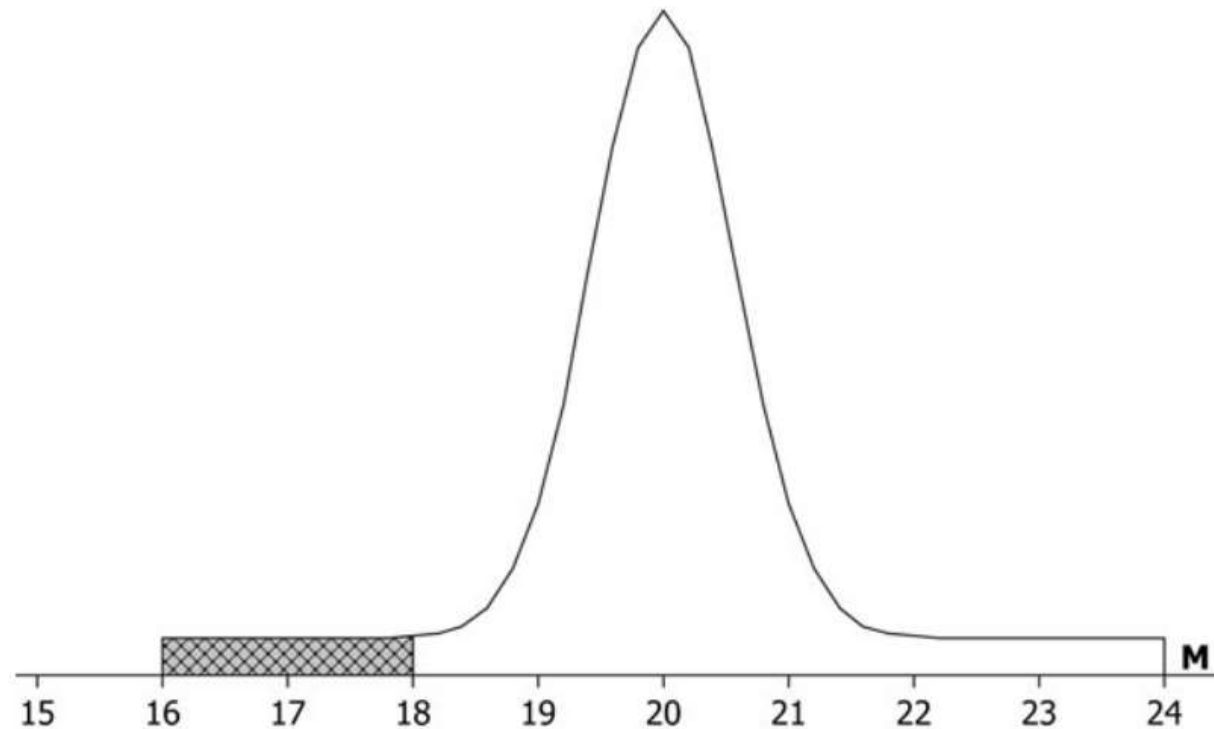
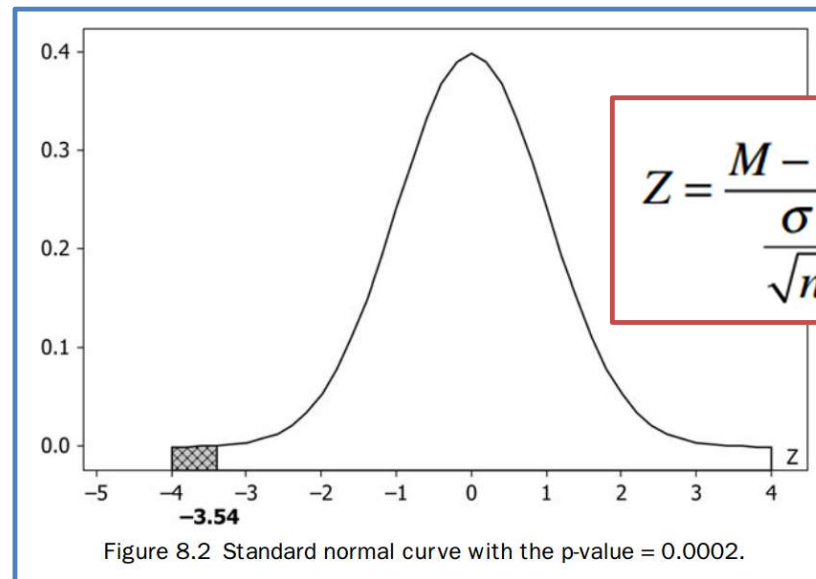
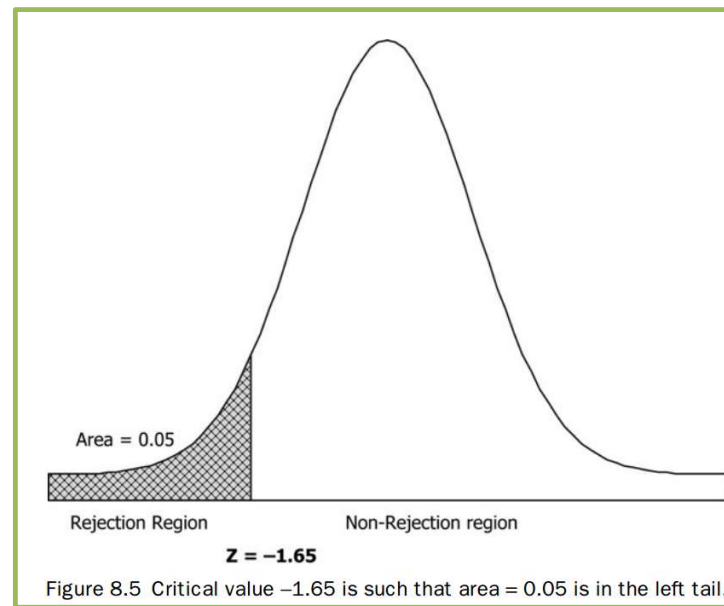


Figure 8.1 Assuming zero affect of the treatment, the left tail area is 0.0002251.



$$Z = \frac{M - 20}{\frac{\sigma}{\sqrt{n}}} = \frac{18 - 20}{\frac{4}{\sqrt{50}}} = -3.54$$



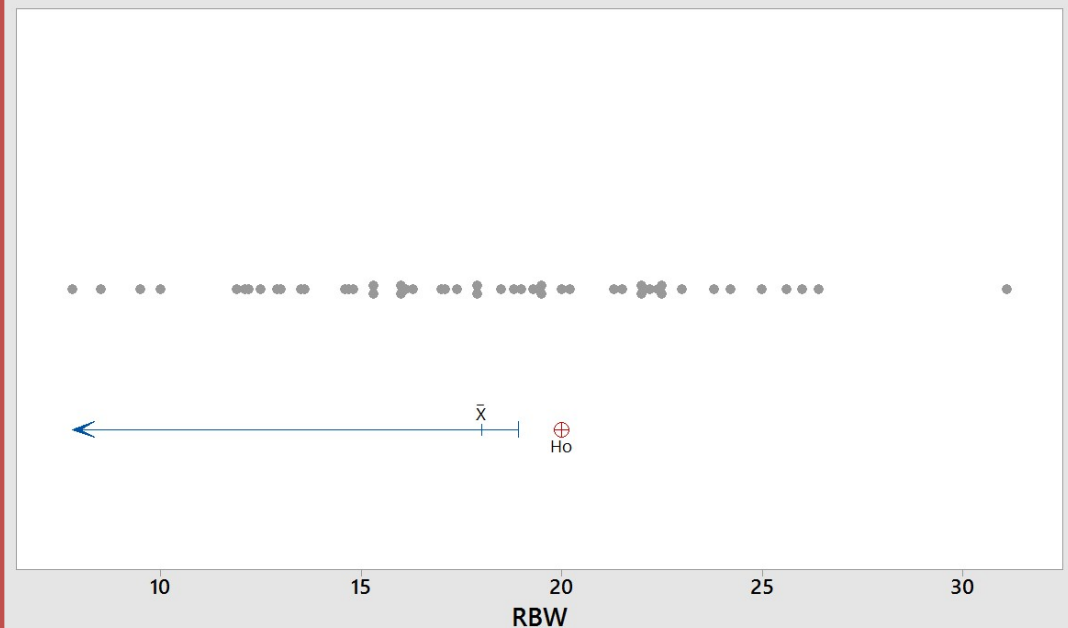
ID	RBW
1	12,2
2	7,8
3	16,1
4	22,4
5	22,2
6	19,5
7	14,7
8	22,0
9	9,5
10	12,1
11	17,9
12	25,0
13	19,0
14	18,8
15	16,0
16	16,3
17	31,1
18	15,3
19	12,9
20	17,1
21	20,0
22	16,0
23	11,9
24	17,0
25	12,5
26	8,5
27	13,5
28	20,2
29	22,5
30	25,6
31	23,8
32	21,5
33	19,5
34	14,6
35	14,8
36	10,0
37	21,3
38	13,0
39	17,9
40	19,3
41	23,0
42	22,0
43	18,5
44	26,0
45	24,2
46	26,4
47	13,6
48	15,3
49	17,4
50	22,5

```
library(readxl)
Dados <- readxl::read_excel("Table 8.2 RawBirthWeight.xls")
BSDA::z.test(x=Dados$RBW, sigma.x=4, mu = 20,
             alternative="less", conf.level=.95)
```

### One-sample z-Test

```
data: Dados$RBW
z = -3.5285, p-value = 0.000209
alternative hypothesis: true mean is less than 20
95 percent confidence interval:
    NA 18.93447
sample estimates:
mean of x
    18.004
```

Individual Value Plot of RBW  
(with Ho and 95% Z-confidence interval for the Mean, and StDev = 4)



# Teste z unilateral

## Valor-p unilateral (greater)

A distribuição da estatura do homem brasileiro de 19 anos de 2016 é normal com desvio-padrão  $\sigma = 7$  cm.

Testar se a média populacional  $\mu$  é menor ou igual a  $\mu_0 = 177$  cm hipotetizada pelo pesquisador.

Quatro participantes desse grupo tiveram suas estaturas medidas, cujos valores em centímetro são: 175, 186, 169, 174.

### Hipóteses

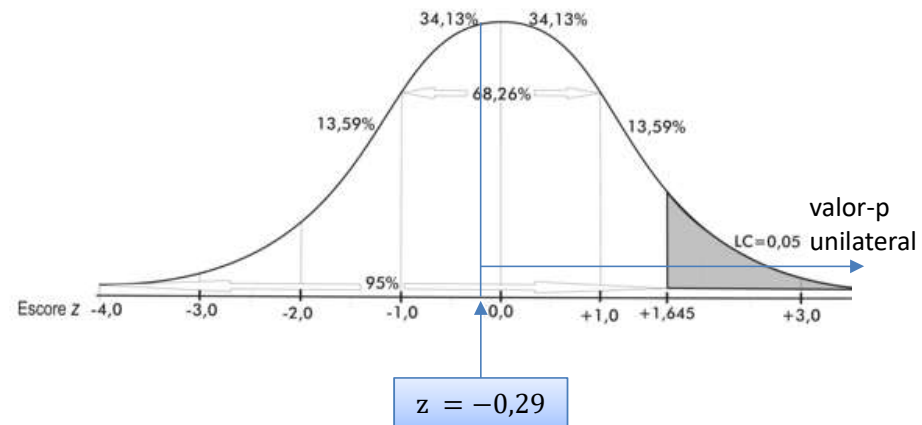
- $H_0: \mu = 177$
- $H_1: \mu > 177$

### Estatísticas

- $\bar{X} = (175+186+169+174)/4 = 176$
- $EP = \frac{\sigma}{\sqrt{n}} = 3,5$
- $IC95(\mu) = [176-1,64 \times 3,5; \infty]$   
=  $[170,24; \infty]$
- Estatística de teste  $z = \frac{\bar{X}-177}{EP} = \frac{176-177}{3,5} = -0,29$

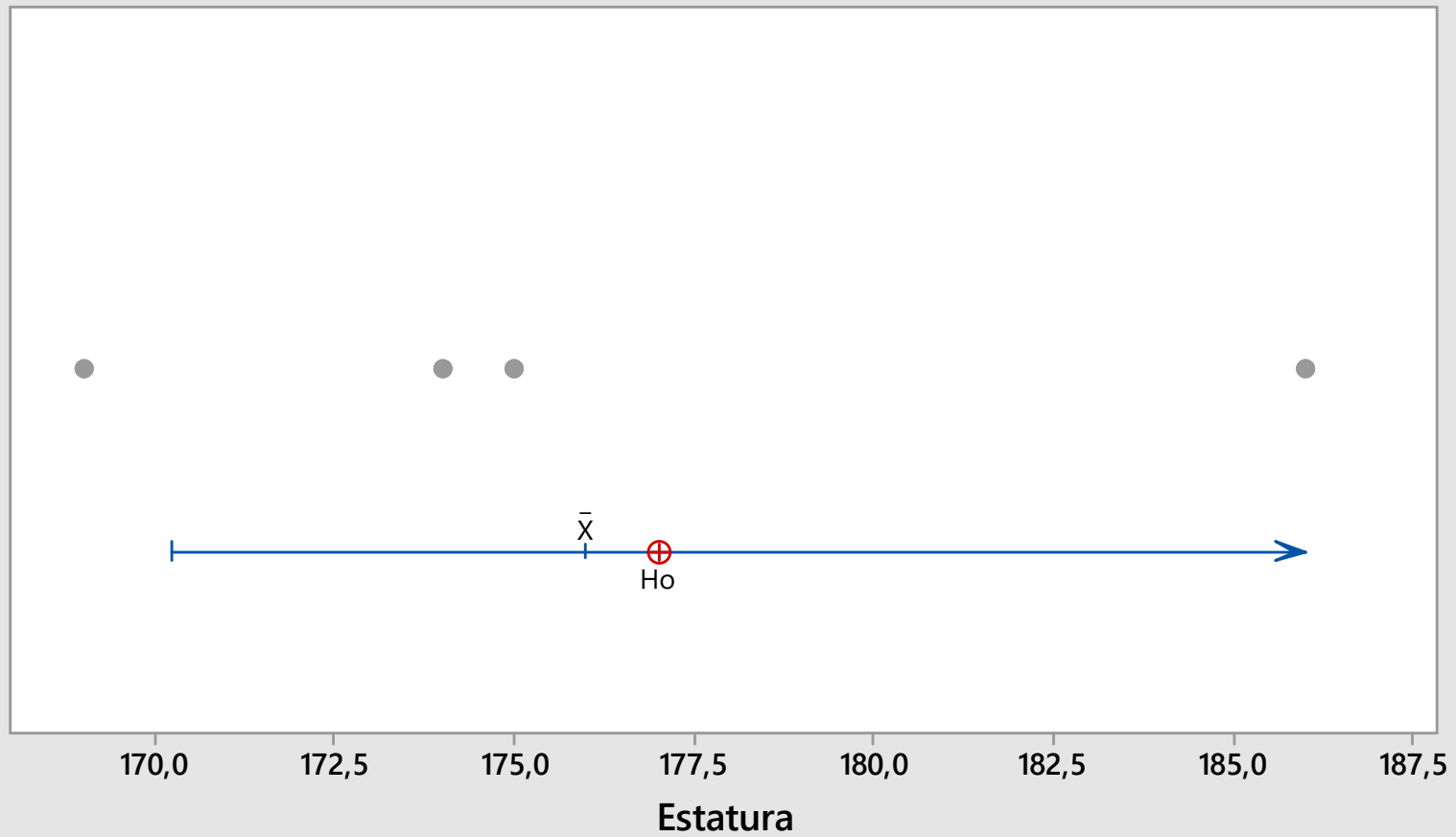
### Decisão

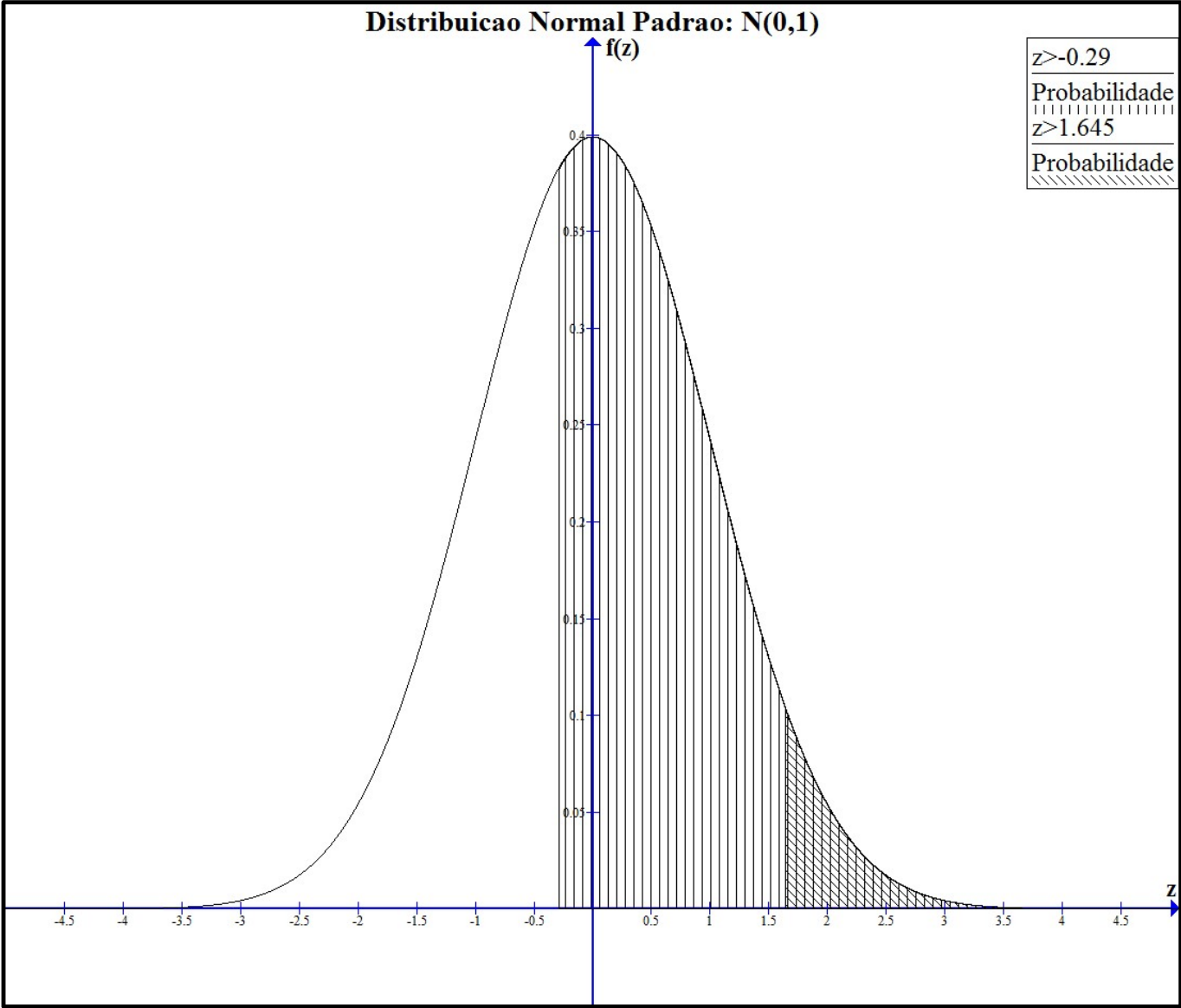
- Como  $z = -0,29 < 1,645$ , não rejeitar  $H_0$  ou
- Como IC95 contém 177, não rejeitar  $H_0$  ou
- Como a probabilidade de escores-z serem mais extremos à direita de -0,29, i.e., o valor-p unilateral = 0,614 (=  $\text{pnorm}(-0.29, \text{mean}=0, \text{sd}=1, \text{lower.tail}=\text{FALSE})$ ) é maior que 5%, não rejeitar  $H_0$
- O valor-p unilateral NÃO é igual à metade do valor-p bilateral = 0,775; é maior que sua metade: 0,388.



# Individual Value Plot of Estatura

(with Ho and 95% Z-confidence interval for the Mean, and StDev = 7)





# Teste z unilateral (greater) para uma condição em R

```
library(BSDA)
estatura <- c(169, 174, 175, 186)
BSDA::z.test(x=estatura, sigma.x=7, mu = 177,
             alternative="greater", conf.level=.95)
```

## One-sample z-Test

```
data:  estatura
z = -0.28571, p-value = 0.6125
alternative hypothesis: true mean is greater than 177
95 percent confidence interval:
 170.243      NA
sample estimates:
mean of x
 176
```

A distribuição da estatura do homem brasileiro de 19 anos de 2016 é normal com desvio-padrão  $\sigma = 7$  cm.

Testar se a média populacional  $\mu$  é menor ou igual a  $\mu_0 = 177$  cm hipotetizada pelo pesquisador.

Quatro participantes desse grupo tiveram suas estaturas medidas, cujos valores em centímetro são: 175, 186, 169, 174.

### *Hipóteses*

- $H_0: \mu = 177$
- $H_1: \mu < 177$

### *Estatísticas*

- $\bar{X} = (175+186+169+174)/4 = 176$
- $EP = \frac{\sigma}{\sqrt{n}} = 3,5$
- $IC95(\mu) = [0; 176+1,64 \times 3,5]$   
= [0; 181,76]
- Estatística de teste  $z = \frac{\bar{X}-177}{EP} = \frac{176-177}{3,5} = -0,29$

### *Decisão*

- Como  $z = -0,29 > -1,645$ , não rejeitar  $H_0$  ou
- Como IC95 contém 177, não rejeitar  $H_0$  ou
- Como a probabilidade de escores-z serem mais extremos à esquerda de -0,29, i.e., o valor-p unilateral = 0,388 (= `pnorm(-0.29, mean=0, sd=1, lower.tail=TRUE)`) é maior que 5%, não rejeitar  $H_0$
- O valor-p unilateral é igual à metade do valor-p bilateral = 0,776.

Teste z unilateral  
Valor-p unilateral (less)



# Teste z unilateral (less) para uma condição em R

```
library(BSDA)
estatura <- c(169, 174, 175, 186)
BSDA::z.test(x=estatura, sigma.x=7, mu = 177,
             alternative="less", conf.level=.95)
```

## One-sample z-Test

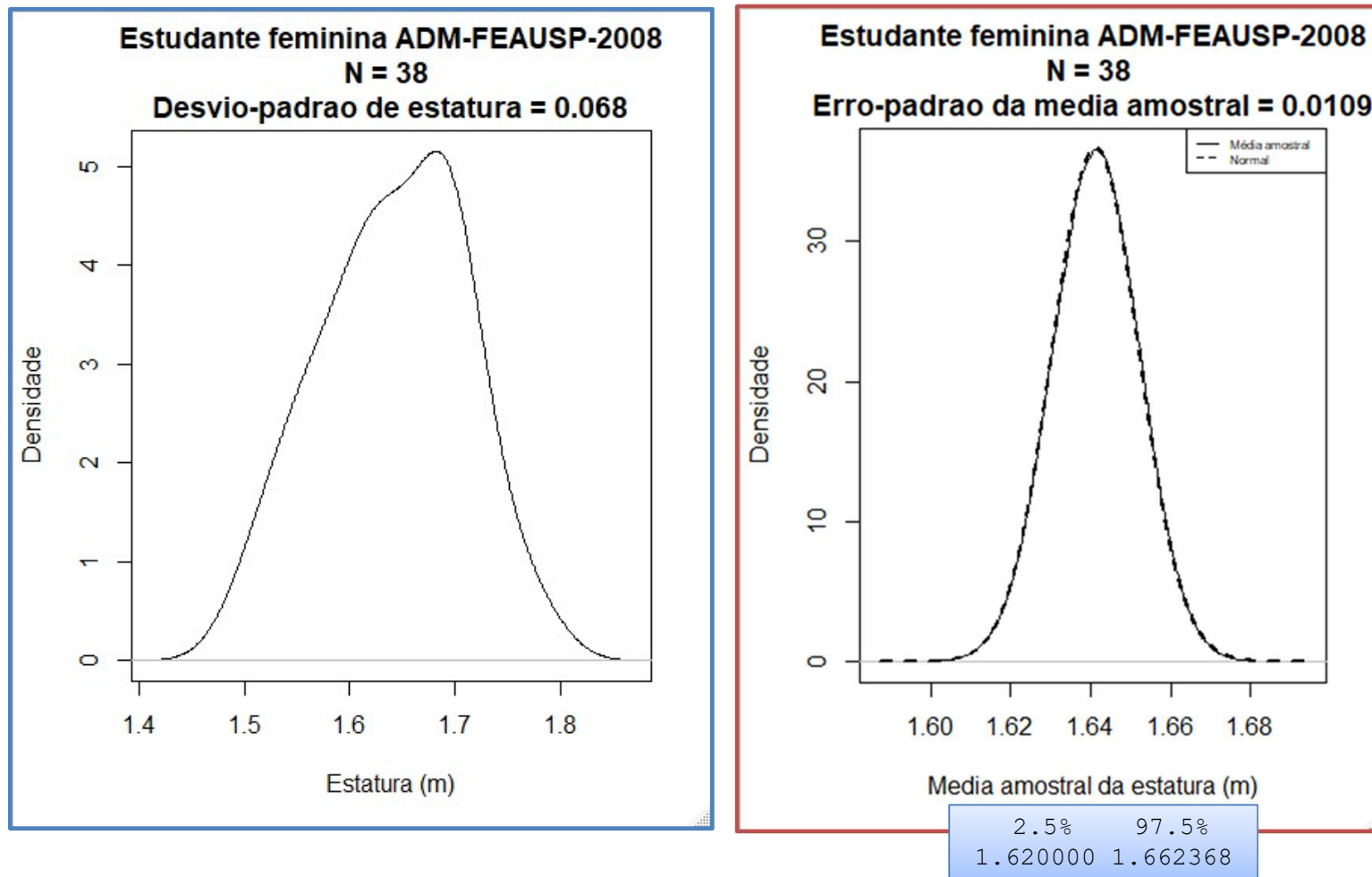
```
data:  estatura
z = -0.28571, p-value = 0.3875
alternative hypothesis: true mean is less than 177
95 percent confidence interval:
    NA 181.757
sample estimates:
mean of x
    176
```

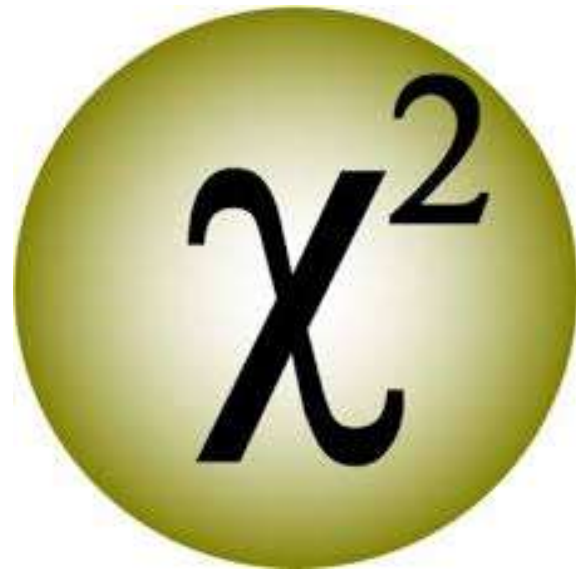
# Reamostragem (*bootstrapping*) da média amostral em R

```
library(readxl)
B <- 1e6; alfa <- 0.05; set.seed(123)
Dados <- readxl::read_excel("Adm2008.xlsx")
Dados <- Dados[, 1:4]
Matriz.Fem <- as.matrix(Dados[Dados$Genero=="Feminino", 3:4])
N.Fem <- nrow(Matriz.Fem)
plot(density(Matriz.Fem[,1], na.rm=TRUE),
     main=paste("Estudante feminina ADM-FEAUSP-2008\nN =", N.Fem,
               "\nDesvio-padroao de estatura =",
               round(sd(Matriz.Fem[,1], na.rm=TRUE), 4)),
     xlab="Estatura (m)", ylab="Densidade")
estat.media.boot.Fem <- replicate(B, mean(sample(Matriz.Fem[,1], replace=TRUE)))
print(mean(Matriz.Fem[,1], na.rm=TRUE))
print(mean(estat.media.boot.Fem, na.rm=TRUE))
quantile(estat.media.boot.Fem, probs=c(alfa/2, 1 - alfa/2))
t.test(Matriz.Fem[,1])$conf.int
plot(density(estat.media.boot.Fem, na.rm=TRUE),
     main=paste("Estudante feminina ADM-FEAUSP-2008\nN =", N.Fem,
               "\nErro-padroao da media amostral =",
               round(sd(estat.media.boot.Fem, na.rm=TRUE), 4)),
     xlab="Media amostral da estatura (m)", ylab="Densidade")
mi <- mean(estat.media.boot.Fem, na.rm=TRUE)
EP <- sd(estat.media.boot.Fem, na.rm=TRUE)
x <- seq(from=mi-5*EP, to=mi+5*EP, by=1e-5)
y <- dnorm(x, mean=mi, sd=EP)
lines(x,y,lwd=2,lty=2)
legend("topright", c("Media amostral","Normal"), lty=1:2, cex=.5)
```

```
[1] 1.641316
> print(mean(Matriz.Fem[,1], na.rm=TRUE))
[1] 1.641316
> print(mean(estat.media.boot.Fem, na.rm=TRUE))
[1] 1.641323
> quantile(estat.media.boot.Fem, probs=c(alfa/2, 1 - alfa/2))
      2.5%      97.5%
1.620000 1.662368
> t.test(Matriz.Fem[,1])$conf.int
[1] 1.618968 1.663663
attr(,"conf.level")
[1] 0.95
```

# Reamostragem (*bootstrapping*) da média amostral em R





# **TESTE QUI-QUADRADO DE DESVIO-PADRÃO POPULACIONAL**

## Teste qui-quadrado bilateral de desvio-padrão populacional para uma condição em R

### *Teste*

- Desvio-padrão  $\sigma = 7$  cm hipotetizado

### *Suposições*

- Estatura tem distribuição normal
- $n = 4$  observações independentes: 169, 174, 175, 186
- Teste bilateral
- Nível de confiança de 95% (ou nível de significância adotado de 5%)

### *Hipóteses*

- $H_0: \sigma = 7$
- $H_1: \sigma \neq 7$

### *Estatísticas*

- $S = 7,16$
- $GL = 4 - 1 = 3$
- Estatística de teste qui-quadrado = 3,14
- $IC95(\sigma) = [4,06; 26,70]$

### *Decisão*

- Como o desvio-padrão populacional hipotetizado está dentro do IC95, não rejeitar  $H_0$
- Como a probabilidade dos valores da estatística de teste qui-quadrado serem mais extremos que 3,14, i.e., o valor-p bilateral = 0,741 é maior que 5%, não rejeitar  $H_0$

# Teste qui-quadrado bilateral de desvio-padrão populacional para uma condição em R

```
library(EnvStats)
estatura <- c(169, 174, 175, 186)
out <- EnvStats::varTest(x=estatura, sigma.squared = 7^2,
  alternative="two.sided", conf.level = 0.95)
print(out)
print(sqrt(out$conf.int))
```

```
Results of Hypothesis Test
-----
Null Hypothesis:                variance = 49
Alternative Hypothesis:         True variance is not equal to 49
Test Name:                      Chi-Squared Test on Variance
Estimated Parameter(s):         variance = 51.33333
Data:                           estatura
Test Statistic:                 Chi-Squared = 3.142857
Test Statistic Parameter:       df = 3
P-value:                        0.7402391
95% Confidence Interval:        LCL = 16.4734
                                UCL = 713.6393

> print(sqrt(out$conf.int))
      LCL      UCL
4.058744 26.714029
attr(,"conf.level")
[1] 0.95
```

# Teste qui-quadrado unilateral (greater) de desvio-padrão populacional para uma condição em R

## *Teste*

- Desvio-padrão  $\sigma = 7$  cm hipotetizado

## *Suposições*

- Estatura tem distribuição normal
- $n = 4$  observações independentes: 169, 174, 175, 186
- Teste bilateral
- Nível de confiança de 95% (ou nível de significância adotado de 5%)

## *Hipóteses*

- $H_0: \sigma = 7$
- $H_1: \sigma > 7$

## *Estatísticas*

- $S^2 = 7,16$
- $GL = 4 - 1 = 3$
- Estatística de teste qui-quadrado = 3,14
- $IC95(\sigma) = [4,44; \infty]$

## *Decisão*

- Como o desvio-padrão populacional hipotetizado está dentro do IC95, não rejeitar  $H_0$
- Como a probabilidade dos valores da estatística de teste qui-quadrado serem mais extremos que 3,14, i.e., o valor-p unilateral à direita = 0,371 é maior que 5%, não rejeitar  $H_0$

# Teste qui-quadrado de desvio-padrão bilateral para uma condição em R

```
library(EnvStats)
estatura <- c(169, 174, 175, 186)
out <- EnvStats::varTest(x=estatura, sigma.squared = 7^2,
                        alternative="greater", conf.level = 0.95)

print(out)
print(sqrt(out$conf.int))
```

```
Results of Hypothesis Test
-----
Null Hypothesis:                variance = 49
Alternative Hypothesis:         True variance is greater than 49
Test Name:                      Chi-Squared Test on Variance
Estimated Parameter(s):        variance = 51.33333
Data:                          estatura
Test Statistic:                 Chi-Squared = 3.142857
Test Statistic Parameter:       df = 3
P-value:                        0.3701195
95% Confidence Interval:        LCL = 19.70638
                                UCL =      Inf

> print(sqrt(out$conf.int))
      LCL      UCL
4.439187      Inf
attr(,"conf.level")
[1] 0.95
```



# Teste qui-quadrado unilateral (less) de desvio-padrão populacional para uma condição em R

## *Teste*

- Desvio-padrão  $\sigma = 7$  cm hipotetizado

## *Suposições*

- Estatura tem distribuição normal
- $n = 4$  observações independentes: 169, 174, 175, 186
- Teste bilateral
- Nível de confiança de 95% (ou nível de significância adotado de 5%)

## *Hipóteses*

- $H_0: \sigma = 7$
- $H_1: \sigma < 7$

## *Estatísticas*

- $S^2 = 7,16$
- $GL = 4 - 1 = 3$
- Estatística de teste qui-quadrado = 3,14
- $IC95(\sigma) = [0;20,92]$

## *Decisão*

- Como o desvio-padrão populacional hipotetizado está dentro do IC95, não rejeitar  $H_0$
- Como a probabilidade dos valores da estatística de teste qui-quadrado serem mais extremos que 3,14, i.e., o valor-p unilateral à direita = 0,63 é maior que 5%, não rejeitar  $H_0$

# Teste qui-quadrado de desvio-padrão bilateral para uma condição em R

```
library(EnvStats)
estatura <- c(169, 174, 175, 186)
out <- EnvStats::varTest(x=estatura, sigma.squared = 7^2,
                        alternative="less", conf.level = 0.95)

print(out)
print(sqrt(out$conf.int))
```

```
Results of Hypothesis Test
-----
Null Hypothesis:                variance = 49
Alternative Hypothesis:         True variance is less than 49
Test Name:                     Chi-Squared Test on Variance
Estimated Parameter(s):        variance = 51.33333
Data:                          estatura
Test Statistic:                Chi-Squared = 3.142857
Test Statistic Parameter:      df = 3
P-value:                       0.6298805
95% Confidence Interval:       LCL = 0.0000
                                UCL = 437.6911

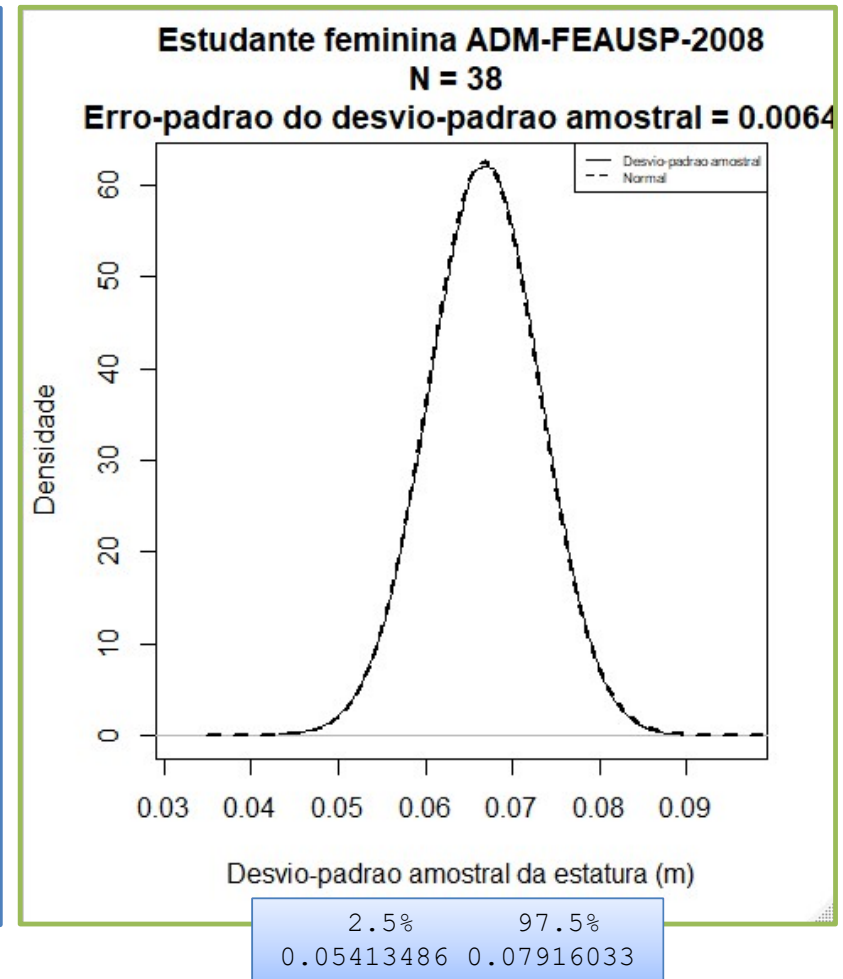
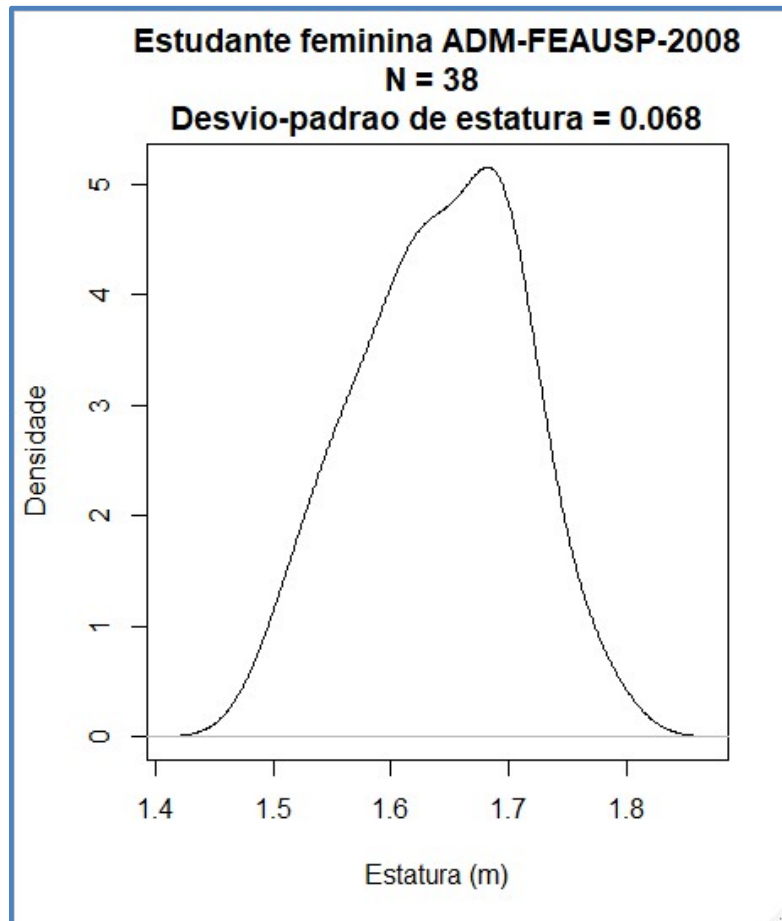
> print(sqrt(out$conf.int))
      LCL      UCL
0.00000 20.92107
attr(,"conf.level")
[1] 0.95
```

# Reamostragem do desvio-padrão amostral em R

```
library(readxl)
library(EnvStats)
B <- 1e6; alfa <- 0.05; set.seed(123)
Dados <- readxl::read_excel("Adm2008.xlsx")
Dados <- Dados[, 1:4]
Matriz.Fem <- as.matrix(Dados[Dados$Genero=="Feminino", 3:4])
N.Fem <- nrow(Matriz.Fem)
plot(density(Matriz.Fem[,1], na.rm=TRUE),
     main=paste("Estudante feminina ADM-FEAUSP-2008\nN =", N.Fem,
               "\nDesvio-padrão de estatura =",
               round(sd(Matriz.Fem[,1], na.rm=TRUE), 4)),
     xlab="Estatura (m)", ylab="Densidade")
estat.media.boot.Fem <- replicate(B, sd(sample(Matriz.Fem[,1], replace=TRUE)))
print(sd(Matriz.Fem[,1], na.rm=TRUE))
print(mean(estat.media.boot.Fem, na.rm=TRUE))
quantile(estat.media.boot.Fem, probs=c(alfa/2, 1 - alfa/2))
ICDP <- EnvStats::varTest(Matriz.Fem[,1])$conf.int
cat("IC95%(DP.Fem) = [", round(sqrt(ICDP[1]),4), ";", round(sqrt(ICDP[2]),4), "]\n")
plot(density(estat.media.boot.Fem, na.rm=TRUE),
     main=paste("Estudante feminina ADM-FEAUSP-2008\nN =", N.Fem,
               "\nErro-padrão do desvio-padrão amostral =",
               round(sd(estat.media.boot.Fem, na.rm=TRUE), 4)),
     xlab="Desvio-padrão amostral da estatura (m)", ylab="Densidade")
mi <- mean(estat.media.boot.Fem, na.rm=TRUE)
EP <- sd(estat.media.boot.Fem, na.rm=TRUE)
x <- seq(from=mi-5*EP, to=mi+5*EP, by=1e-5)
y <- dnorm(x, mean=mi, sd=EP)
lines(x,y,lwd=2,lty=2)
legend("topright", c("Desvio-padrão amostral", "Normal"), lty=1:2, cex=.5)
```

```
> print(sd(Matriz.Fem[,1], na.rm=TR
[1] 0.06798931
> print(mean(estat.media.boot.Fem,
[1] 0.0667759
> quantile(estat.media.boot.Fem, pr
          2.5%      97.5%
0.05413486 0.07916033
> ICDP <- EnvStats::varTest(Matriz.
> cat("IC95%(DP.Fem) = [", round(sq
IC95%(DP.Fem) = [ 0.0554 ; 0.088 ]
```

# Reamostragem do desvio-padrão amostral em R



**Tabela 4.1**

Tipo de análise da potência	O que se pretende determinar	O que se deve especificar
Análise da potência <i>a priori</i>	$n$	$\alpha$ , $1 - \beta$ e ES
Análise da potência <i>post hoc</i>	$1 - \beta$	$\alpha$ , ES e $n$
Análise da potência de compromisso	$\alpha$ e $1 - \beta$	ES, $n$ e $\frac{\beta}{\alpha}$
Análise de sensibilidade	ES mínimo	$\alpha$ , $1 - \beta$ e $n$
→ Análise do critério	$\alpha$	$1 - \beta$ , ES e $n$

COELHO, JP *et al.* (2008) *Inferência Estatística: com utilização do SPSS e G\*Power*. Lisboa: Sílabo.

---

**RESEARCH ARTICLE**

# A Practical Primer To Power Analysis for Simple Experimental Designs

Marco Perugini, Marcello Gallucci and Giulio Costantini

---

Power analysis is an important tool to use when planning studies. This contribution aims to remind readers what power analysis is, emphasize why it matters, and articulate when and how it should be used. The focus is on applications of power analysis for experimental designs often encountered in psychology, starting from simple two-group independent and paired groups and moving to one-way analysis of variance, factorial designs, contrast analysis, trend analysis, regression analysis, analysis of covariance, and mediation analysis. Special attention is given to the application of power analysis to moderation designs, considering both dichotomous and continuous predictors and moderators. Illustrative practical examples based on G\*Power and R packages are provided throughout the article. Annotated code for the examples with R and dedicated computational tools are made freely available at a dedicated web page (<https://github.com/mcfanda/primerPowerIRSP>). Applications of power analysis for more complex designs are briefly mentioned, and some important general issues related to power analysis are discussed.

---

**Keywords:** power analysis; effect size; moderation; sensitivity analysis; uncertainty

---



# LIMITS OF RETROSPECTIVE POWER ANALYSIS

PATRICK D. GERARD,<sup>1</sup> Experimental Statistics Unit, Box 9653, Mississippi State University, Mississippi State, MS 39762, USA  
 DAVID R. SMITH, U.S. Geological Service, Biological Resources Division, Leetown Science Center, Kearneysville, WV 25430, USA  
 GOVINDA WEERAKKODY, Department of Mathematics and Statistics, Box 9715, Mississippi State University, Mississippi State, MS 39762, USA

**Abstract:** Power analysis after study completion has been suggested to interpret study results. We present 3 methods of estimating power and discuss their limitations. We use simulation studies to show that estimated power can be biased, extremely variable, and severely bounded. We endorse the practice of computing power to detect a biologically meaningful difference as a tool for study planning but suggest that calculation of confidence intervals on the parameter of interest is the appropriate way to gauge the strength and biological meaning of study results.

1998 *JOURNAL OF WILDLIFE MANAGEMENT* 62(2):801–807

**Key words:** interpretation of results, noncentrality parameter, statistical power, study design.

Table 1. Expected power when power is estimated with the observed test statistic based on simulated comparisons between 3 populations where sample size was 20/population. Six values of power (0.05, 0.11, 0.34, 0.66, 0.90, 0.98) were created by setting  $\alpha = 0.05$  and altering population means. Population means were distributed normally with unit variance. Three estimators of power were calculated for each replication; estimators are discussed in detail in the text. Expected power,  $E(P)$ , and expected upper and lower bounds for the 95% confidence interval (CI) were computed by averaging over 500 replications. For each replicate and the simple plug-in ( $P_p$ ) and corrected ( $P_{bc}$ ) estimators, the bounds for the 95% CI were the appropriate percentiles from a bootstrap sample of size 400. For the percentile estimator ( $P_{mi}$ ), the 95% CI was based on exact methods.

True power	Simple plug-in ( $P_p$ )		Corrected ( $P_{bc}$ )		Percentile ( $P_{mi}$ )	
	$E(P_p)$	95% CI	$E(P_{bc})$	95% CI	$E(P_{mi})$	95% CI
0.05	0.236	(0.072, 0.870)	0.126	(0.052, 0.790)	0.163	(0.051, 0.671)
0.11	0.271	(0.081, 0.890)	0.151	(0.055, 0.820)	0.194	(0.055, 0.717)
0.34	0.455	(0.133, 0.948)	0.309	(0.073, 0.912)	0.375	(0.074, 0.870)
0.66	0.662	(0.259, 0.981)	0.538	(0.153, 0.966)	0.601	(0.148, 0.950)
0.90	0.844	(0.462, 0.995)	0.770	(0.326, 0.991)	0.811	(0.309, 0.987)
0.98	0.948	(0.678, 0.999)	0.915	(0.561, 0.998)	0.934	(0.533, 0.998)

# Análise de poder retrospectivo (plug in) de ANOVA unifatorial independente balanceada (Gerard *et al.* (1998)

Análise de poder retrospectivo de ANOVA unifatorial independente balanceada.R

```
library(MBESS)
k <- 3
n <- 20
alfa <- 0.05
dfn <- k - 1
dfd <- k*(n - 1)
Fcrt <- qf(1-alfa, dfn, dfd)
Fobs <- Fcrt*0.99
eta2 <- dfn*Fobs/(dfn*Fobs+dfd)
eta2lims <- MBESS::ci.pvaf(Fobs, dfn, dfd, k*n, 1-alfa)
f2 <- eta2/(1-eta2)
f2.ll <- eta2lims$Lower.Limit.Proportion.of.Variance.Accounted.for/
  (1-eta2lims$Lower.Limit.Proportion.of.Variance.Accounted.for)
f2.ul <- eta2lims$Upper.Limit.Proportion.of.Variance.Accounted.for/
  (1-eta2lims$Upper.Limit.Proportion.of.Variance.Accounted.for)
ncp.p <- dfd*f2 # ou dfn*Fobs
ncp.p.ll <- dfd*f2.ll
ncp.p.ul <- dfd*f2.ul
cat(paste("N =", k*n, "\tFcrt =", round(Fcrt,2), "\tFobs =", round(Fobs,2), "\n"))
poder.p <- 1-pf(Fcrt,dfn, dfd, ncp.p)
cat(paste("\tPoder.p =", round(poder.p,3), "\n"))
poder.p.ll <- 1-pf(Fcrt,dfn, dfd, ncp.p.ll)
cat(paste("\tPoder.p.ll =", round(poder.p.ll,3), "\n"))
poder.p.ul <- 1-pf(Fcrt,dfn, dfd, ncp.p.ul)
cat(paste("\tPoder.p.ul =", round(poder.p.ul,3), "\n"))
sink()
```

```
N = 60          Fcrt = 3.16    Fobs = 3.13
                Poder.p = 0.579
                Poder.p.ll = 0.05
                Poder.p.ul = 0.967
```



DEBATE

Open Access

# Current sample size conventions: Flaws, harms, and alternatives

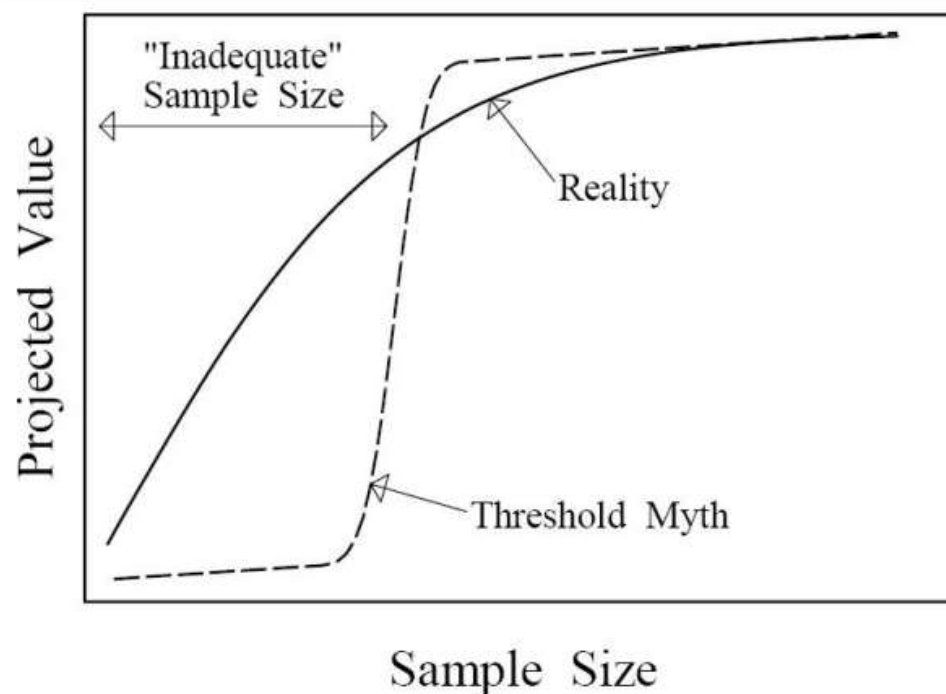
Peter Bacchetti

## Abstract

**Background:** The belief remains widespread that medical research studies must have statistical power of at least 80% in order to be scientifically sound, and peer reviewers often question whether power is high enough.

**Discussion:** This requirement and the methods for meeting it have severe flaws. Notably, the true nature of how sample size influences a study's projected scientific or practical value precludes any meaningful blanket designation of <80% power as "inadequate". In addition, standard calculations are inherently unreliable, and focusing only on power neglects a completed study's most important results: estimates and confidence intervals. Current conventions harm the research process in many ways: promoting misinterpretation of completed studies, eroding scientific integrity, giving reviewers arbitrary power, inhibiting innovation, perverting ethical standards, wasting effort, and wasting money. Medical research would benefit from alternative approaches, including established *value of information* methods, simple choices based on cost or feasibility that have recently been justified, sensitivity analyses that examine a meaningful array of possible findings, and following previous analogous studies. To promote more rational approaches, research training should cover the issues presented here, peer reviewers should be extremely careful before raising issues of "inadequate" sample size, and reports of completed studies should not discuss power.

**Summary:** Common conventions and expectations concerning sample size are deeply flawed, cause serious harm to the research process, and should be replaced by more rational alternatives.



**Figure 1 Qualitative depiction of how sample size influences a study's projected scientific and/or practical value.** A threshold shaped relationship (dashed line) would create a meaningful distinction between adequate and inadequate sample sizes, but such a relation does not exist. The reality (solid line) is qualitatively different, exhibiting diminishing marginal returns. Under the threshold myth, cutting a sample size in half could easily change a valuable study into an inadequate one, but in reality such a cut will always preserve *more* than half of the projected value.

# Referências

- COELHO, JP et al. (2008) Inferência Estatística: com utilização do SPSS e G\*Power. Lisboa: Sílabo.
- DANCEY, C & REIDY, J (2019)  
*Estatística sem Matemática para Psicologia*.  
7ª ed. Porto Alegre: Penso.
- STEPHENS, LJ (2009) *Statistics in Psychology*. NY: McGraw-Hill, Schaum's Outline Series.

