

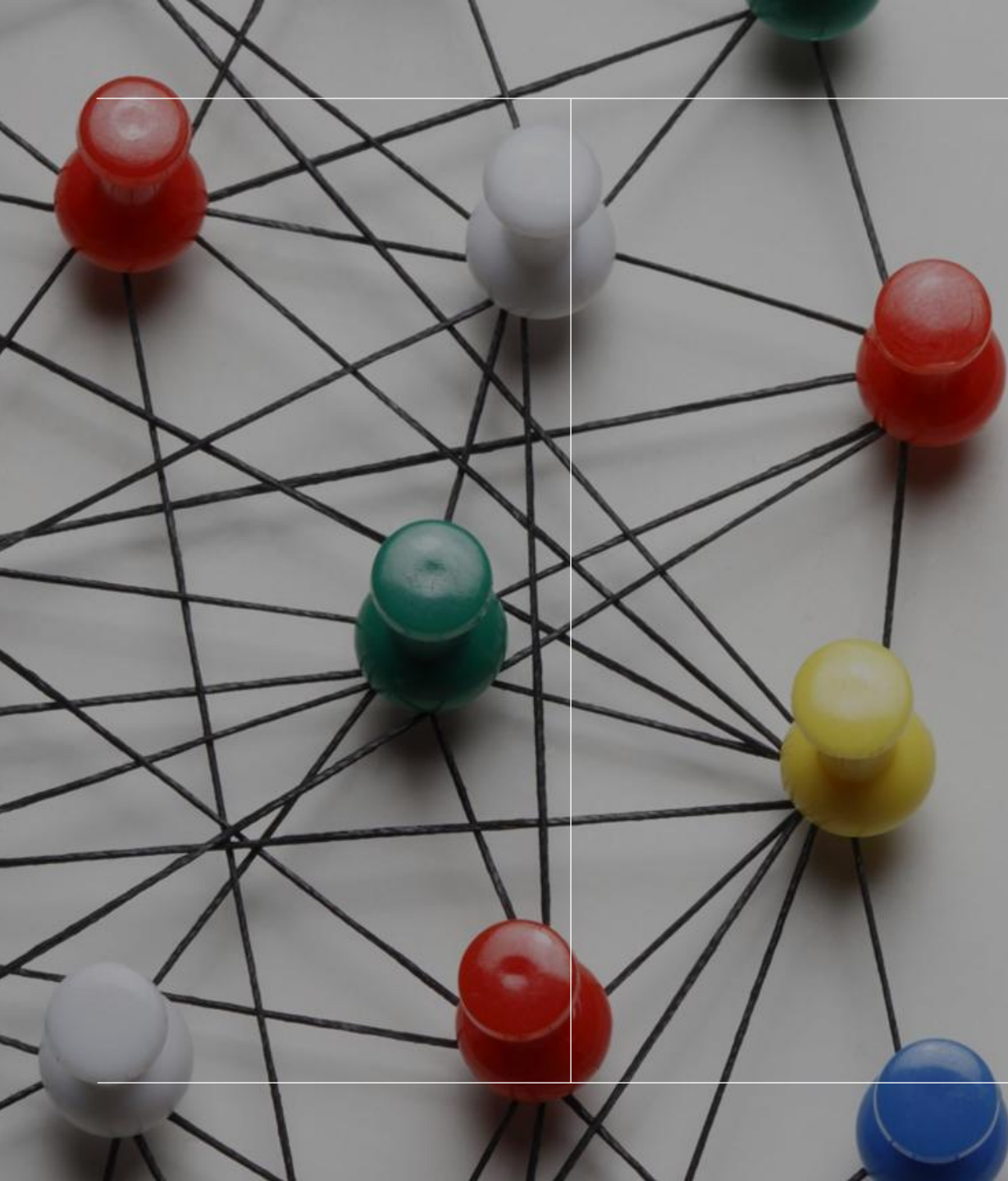


Aulas 6, 7 e 8 Inferência Estatística

PROF. ME. NATAN KLEIN
IA & ANÁLISE DE DADOS
CIÊNCIA DA COMPUTAÇÃO
ATITUS EDUCAÇÃO

Conteúdos

- Inferência
 - Amostras e populações
 - Representatividade
 - Amostragem
 - Distribuições probabilísticas
 - Distribuição normal
 - Desvios da normalidade
 - Distribuição amostral
 - Teorema do limite central
 - Erro padrão
 - Intervalos de confiança
 - Testagem de hipóteses
 - Bootstrap
-



Inferência

Tipos de raciocínio

- Existem três tipos básicos de raciocínio:

Característica	Raciocínio Dedutivo	Raciocínio Indutivo	Raciocínio Abduativo
Ponto de Partida	Premissas gerais	Observações específicas	Observação de um fato surpreendente ou incompleto
Conclusão	Necessariamente verdadeira se as premissas forem verdadeiras	Provavelmente verdadeira, generalização das observações	Melhor explicação possível para o fato observado (hipótese)
Direção do Raciocínio	Do geral para o específico	Do específico para o geral	Do efeito para a possível causa
Certeza da Conclusão	Conclusiva	Probabilística	Plausível, mas não garantida
Objetivo Principal	Provar uma conclusão com base em fatos conhecidos	Descobrir padrões e formular generalizações	Explicar um fenômeno ou encontrar a melhor hipótese
Exemplo Simplificado	Todo A é B. C é A. Logo, C é B.	Cisne 1 é branco, cisne 2 é branco... Logo, todos os cisnes são brancos.	A grama está molhada. Pode ter chovido.

Indução, inferência e hipóteses

- A inferência é um exemplo de raciocínio indutivo ... as conclusões que inferimos são, sempre, provavelmente verdadeiras.
- Uma hipótese é uma inferência que fazemos para explicar algo.
- Podemos usar estatística para formular hipóteses (análise exploratória de dados) ou testar hipóteses (estatística inferencial).



Abordagens científicas

Abordagem confirmatória,
guiada pela
teoria.

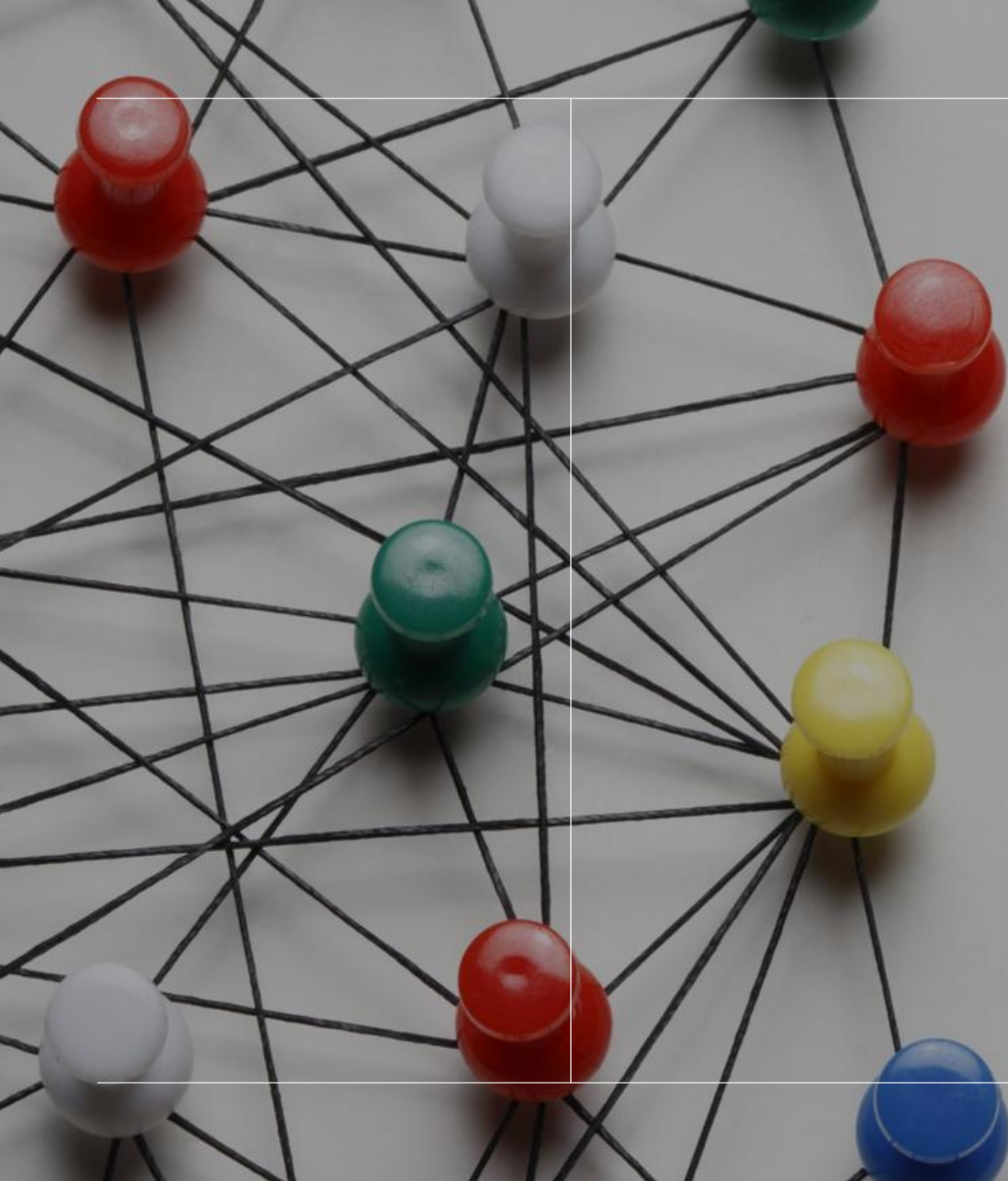


Abordagem exploratória,
guiada pelos
dados.



Vamos praticar!





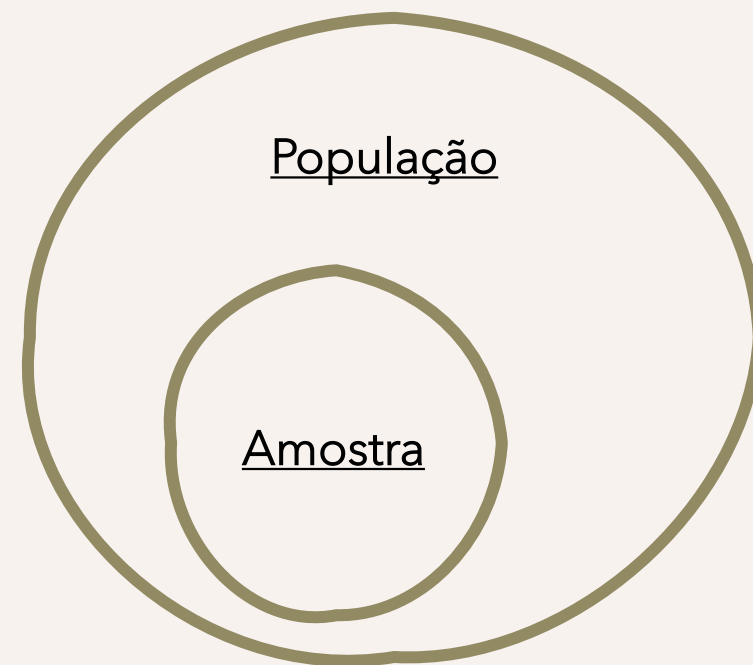
Amostras e populações

Exemplos

População	Amostra
População Brasileira	Subconjunto de residentes do Brasil
Alunos de uma faculdade	% de alunos de cada curso
Sangue de uma pessoa	Amostra de sangue em jejum
Tempo	Períodos aleatórios
Linha de produção (xícaras)	Peças em períodos aleatórios de tempo
Características psicológicas (nível de disciplina)	Comportamentos específicos

Definições

- População: um grupo completo de referência sobre o qual queremos tirar conclusões.
- Amostra: subconjunto da população que usamos para fazer inferências.



Parâmetros

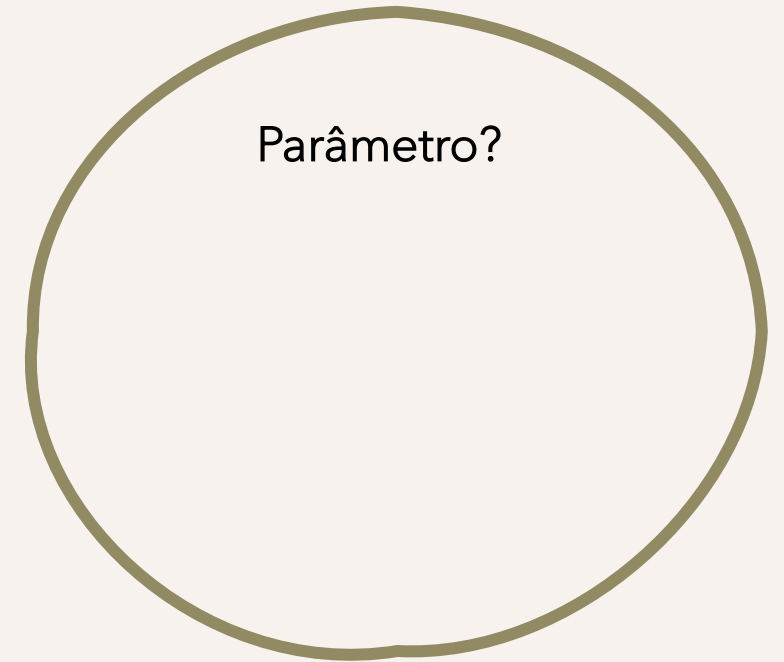
- Queremos calcular métricas que digam algo sobre uma característica da população (variável).
- Chamamos essas métricas da população de parâmetros.
- Exemplo de parâmetro: consumo médio de refrigerantes por dia dos brasileiros.



Média da pop: μ

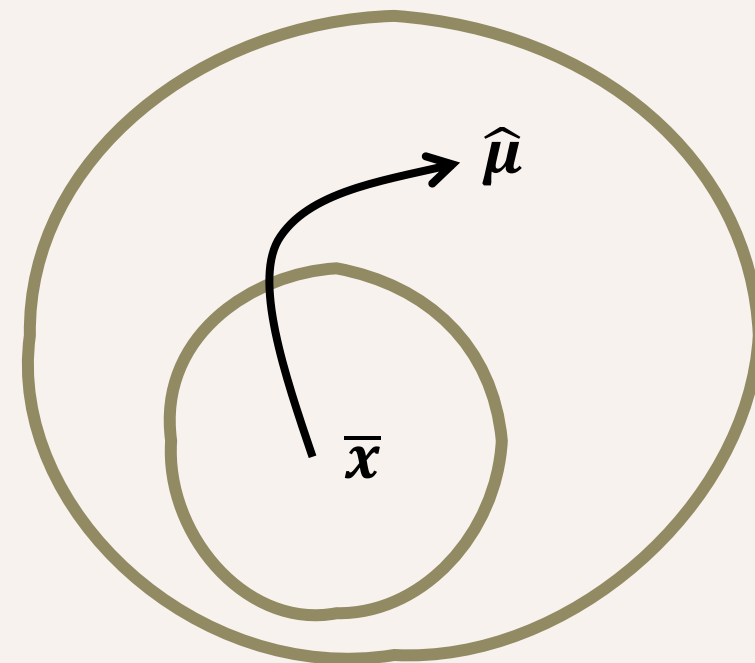
Parâmetros

- Mas não temos acesso direto a muitos parâmetros:
 - Nível de depressão de estudantes de graduação
 - Taxa de suicídios por pessoa no Brasil
 - Nível de felicidade das pessoas de um país
 - Número de erros em uma linha de produção ao longo de um ano
 - ...



Estimativa de parâmetros

- O que fazemos para descobrir um valor aproximado para o parâmetro é estimá-lo ou inferi-lo a partir de uma amostra que represente a população!
- Nessa amostra, calculamos uma estatística e usamos ela como estimativa do parâmetro (melhor chute possível).



Exemplo de estimativa de parâmetros

- Exemplo: Pesquisa Nacional de Saúde
-

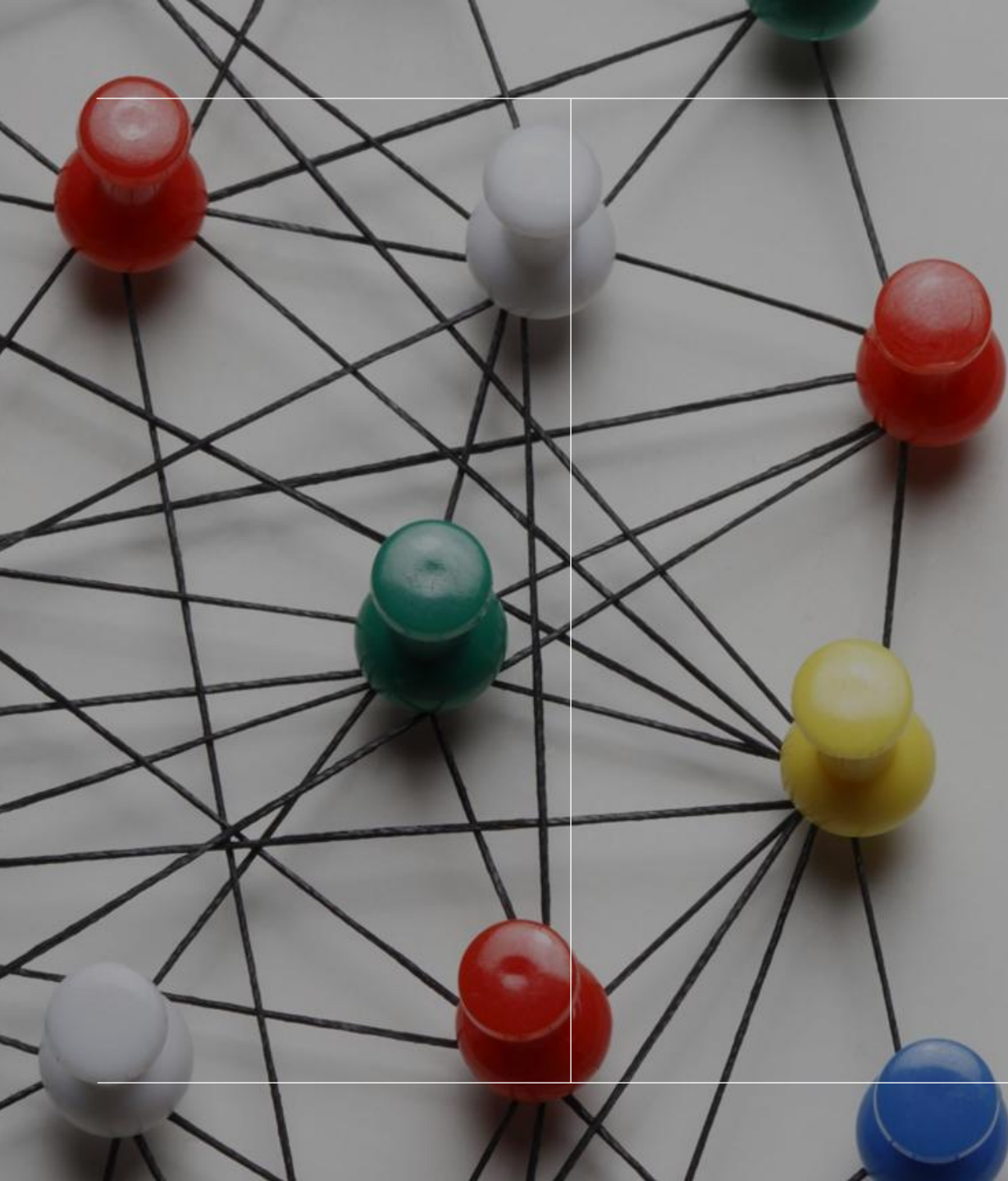
Convenções

- Existem muitos tipos de parâmetros e estatísticas. Uma convenção é usar:
 - Parâmetros \rightarrow letras gregas: $\alpha \beta \gamma \delta \varepsilon \epsilon \zeta \eta \theta \vartheta \iota \kappa \lambda \mu \nu \xi \omicron \pi \varpi \rho \varphi \phi \chi \psi \omega \dots$
 - Estatísticas \rightarrow letras latinas: $a b c d e f g h i j k l \dots$
 - Existem estatísticas/parâmetros:
 - **Univariados:** calculados com apenas uma variável.
 - **Bivariados:** calculados com duas variáveis.
 - **Multivariados:** calculados com mais de duas variáveis.
-

Exemplos de parâmetros/estatísticas

- Ver tabela!

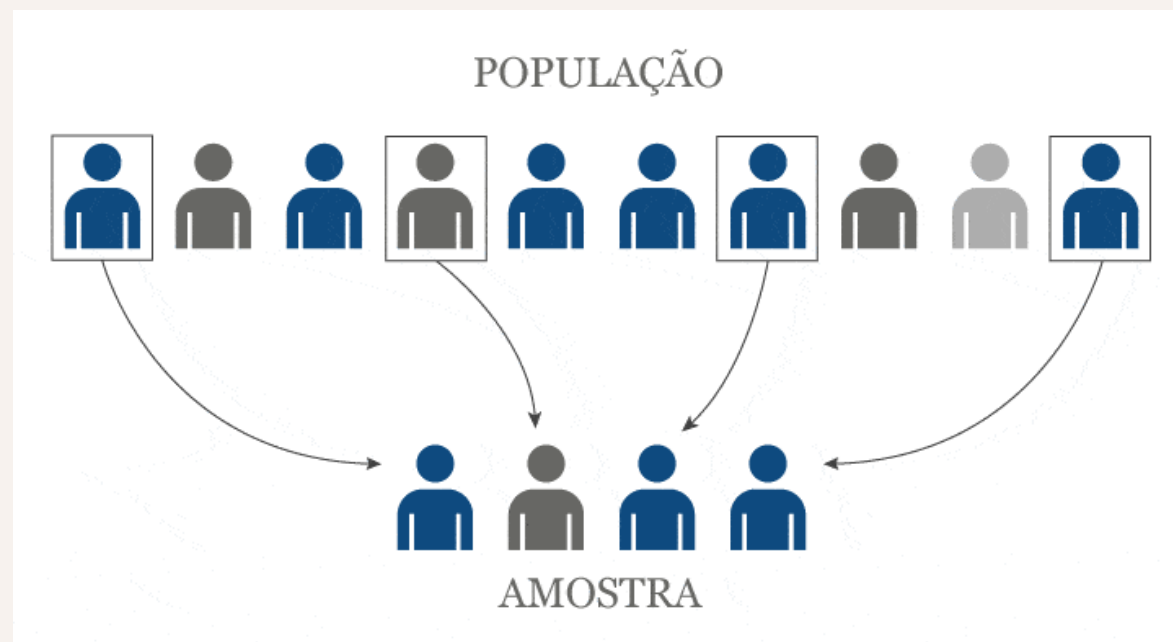




Representatividade

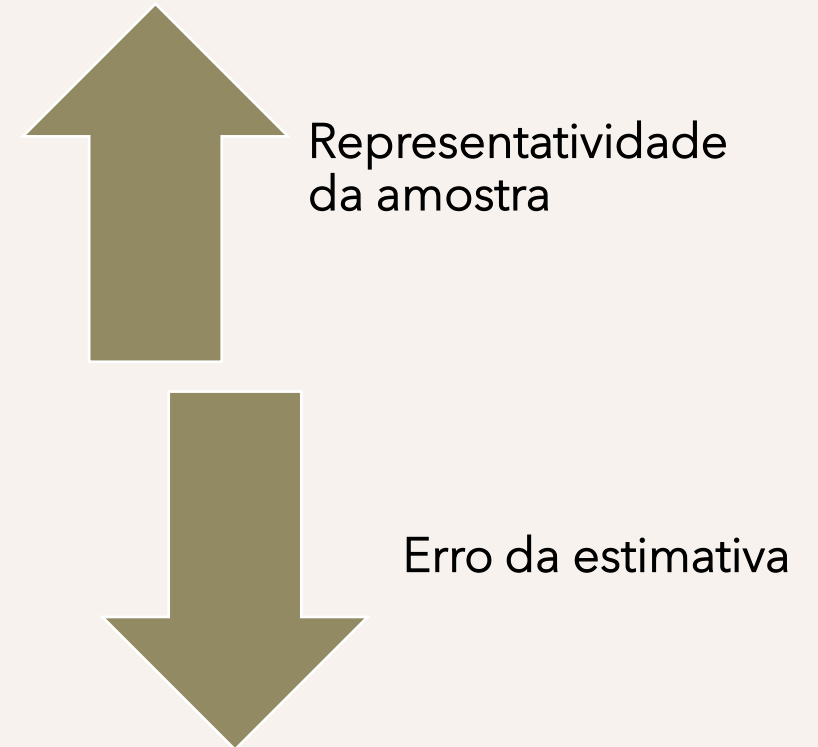
Definição

- Queremos extrair uma amostra que seja representativa da população.
- Uma amostra representativa da população é aquela que reproduz as características fundamentais da população e nos permite calcular estatísticas que estimam parâmetros com o mínimo de erro.



Confiança

- Como posso saber que minhas estatísticas amostrais são medidas confiáveis, fidedignas, dos parâmetros da população?
- As estimativas serão tão boas quanto a amostra for representativa da população, i.e., amostra representativa = boas estimativas.



Erros e vieses

- O erro de estimativa de um parâmetro pode ser expresso pela equação:

$$\textit{erro} = \textit{parâmetro} - \textit{estatística}$$

- Uma amostra não-representativa é fundamentalmente diferente da população em algum aspecto importante, o que enviesa (introduz erro) nas estimativas.
-

Erros e vieses

Tipo de Viés	Definição	Exemplo
Viés de seleção	Ocorre quando a amostra não é selecionada aleatoriamente, favorecendo certos grupos.	Pesquisa sobre hábitos alimentares em academias, favorecendo pessoas preocupadas com a saúde.
Viés de não resposta	Ocorre quando certos grupos têm maior probabilidade de não responder à pesquisa.	Pesquisa online sobre privacidade com baixa resposta de idosos, menos familiarizados com a internet.
Viés de sobrevivência	Amostra composta apenas por "sobreviventes" de um processo, ignorando os que não sobreviveram.	Estudo de sucesso de empresas iniciantes, incluindo apenas as que ainda operam.
Viés de conveniência	Amostra selecionada pela facilidade de acesso aos participantes.	Pesquisa em shopping center, representando apenas frequentadores de shoppings.
Viés de subcobertura	Alguns membros da população são representados inadequadamente na amostra.	Pesquisa por telefone excluindo pessoas sem telefone fixo ou que não atendem desconhecidos.
Viés de resposta voluntária	Amostra composta apenas por voluntários, geralmente com opiniões fortes.	Pesquisa online sobre produto controverso, atraindo pessoas com opiniões extremas.

Garantias da representatividade

- Existem duas formas de garantir a representatividade de uma amostra:
 - Tamanho da amostra
 - Método de seleção da amostra (amostragem) → **mais importante!**

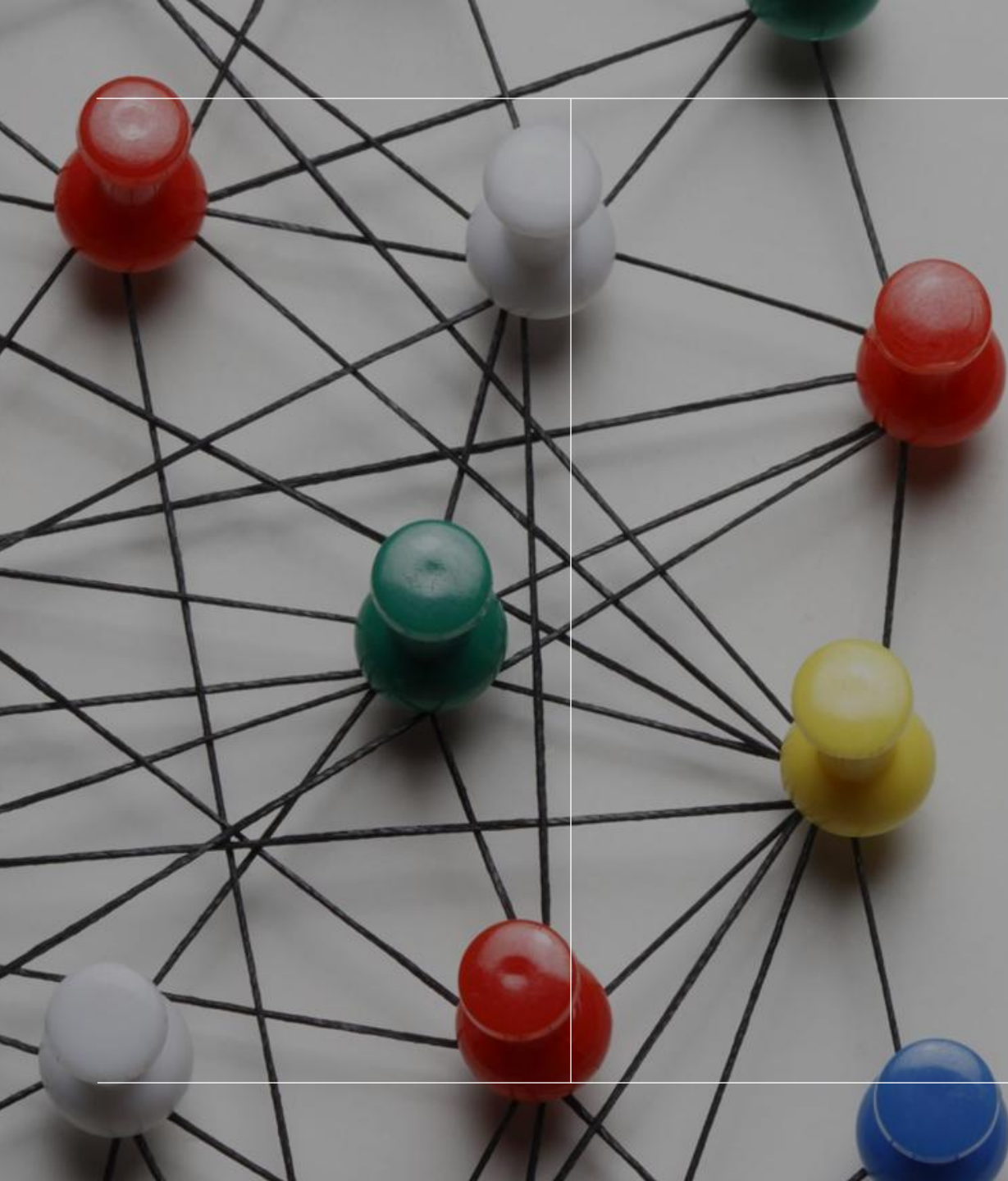
Tamanho da amostra

- Estatísticas calculadas em amostras muito pequenas dificilmente serão representativas da população.
 - Não há um critério para o tamanho mínimo de uma amostra porque isso varia de um contexto para outro, mas $n > 30$ é aconselhável.
 - Existem formas de descobrir o tamanho ideal de uma amostra por meio de cálculos de tamanho amostral (estatísticas de poder amostral).
-

Tamanho da amostra

- Vamos usar o banco de dados dos passageiros do Titanic para calcular a média de idade dos passageiros com diferentes tamanhos de amostra selecionadas aleatoriamente!

Amostra	Média	Erro
Total (N)	29,70	-
$n = 10$	37,28	+7,58
$n = 30$	34,73	+5,03
$n = 60$	31,44	+1,74
$n = 100$	30,38	+0,68



Amostragem

Amostragem aleatória simples

- O melhor método para selecionar objetos da população para nossa amostra é a seleção aleatória:
 - Garante a máxima representatividade.
 - Cada elemento da população é conhecido e é selecionado aleatoriamente para entrar na amostra.
 - Todos tem a mesma chance de entrar na amostra.
-

Amostragem aleatória simples

- Problemas:

- Necessita uma base amostral com todos os elementos da população.
- Processo custoso.
- Vieses ainda podem atrapalhar o processo.

- Exemplos:

- Estudantes da Atitus
 - Pacientes de um hospital
 - População geral
-

Literary Digest vs Gallup

REPRODUCED FROM HOLDINGS AT THE FRANKLIN D. ROOSEVELT LIBRARY

The Literary Digest

NEW YORK

OCTOBER 31, 1936

Topics of the day

LONDON, 1,293,669; ROOSEVELT, 972,897

Final Returns in The Digest's Poll of Ten Million Voters

We all, the great battle of the ballots in the Poll of ten million voters, scattered throughout the forty-eight States of the Union, is now finished, and in the table below we record the figures received up to the hour of going to press.

These figures are exactly as received from more than one in every five voters polled in our country—they are neither weighted, adjusted nor interpreted.

Never before in an experience covering more than a quarter of a century in taking polls have we received so many different varieties of criticism—praise from many; condemnation from many others—and yet it has been just of the same type that has come to us every time a Poll has been taken in all these years.

A telegram from a newspaper in California asks: "Is it true that Mr. Hearst has purchased The Literary Digest?" A telephone message only the day before these lines were written: "Has the Repub-

lican National Committee purchased The Literary Digest?" And all types and varieties, including: "Have the Jews purchased The Literary Digest?" "Is the Pope of Rome a stockholder of The Literary Digest?" And so it goes—all equally absurd and amusing. We could add more to this list, and yet all of these questions in recent days are but repetitions of what we have been experiencing all down the years from the very first Poll.

Problem—Now, are the figures in this Poll correct? In answer to this question we will simply refer to a telegram we sent to a young man in Massachusetts the other day in answer to his challenge to us to wager \$100,000 on the accuracy of our Poll. We wired him as follows:

"For nearly a quarter century, we have been taking Polls of the voters in the forty-eight States, and especially in Presidential years, and we have always merely mailed the ballots, counted and recorded those

returned and let the people of the Nation draw their conclusions as to our accuracy. So far, we have been right in every Poll. Will we be right in the current Poll? That, as Mrs. Roosevelt said concerning the President's reelection, is in the 'lap of the gods.' We never make any claims before election but we respectfully refer you to the opinion of one of the most quoted citizens to-day, the Hon. James A. Farley, Chairman of the Democratic National Committee. This is what Mr. Farley said October 14, 1936:

"Any sane person can not escape the implication of such a gigantic sampling of popular opinion as is embraced in The Literary Digest straw vote. I consider this conclusive evidence as to the desire of the people of this country for a change in the National Government. The Literary Digest poll is an achievement of no little magnitude. It is a Poll fairly and correctly conducted."

In studying the table of the voters from

the statistics and the material in this article are the property of Funk & Wagnalls Company and have been copyrighted by it; neither this whole nor any part thereof may be reprinted or published without the special permission of the copyright owner.

Final Report "Literary Digest" 1936 Presidential Poll

How the Same Voters Voted in the 1932 Election

Roosevelt

How the Same Voters Voted in the 1932 Election

Landon

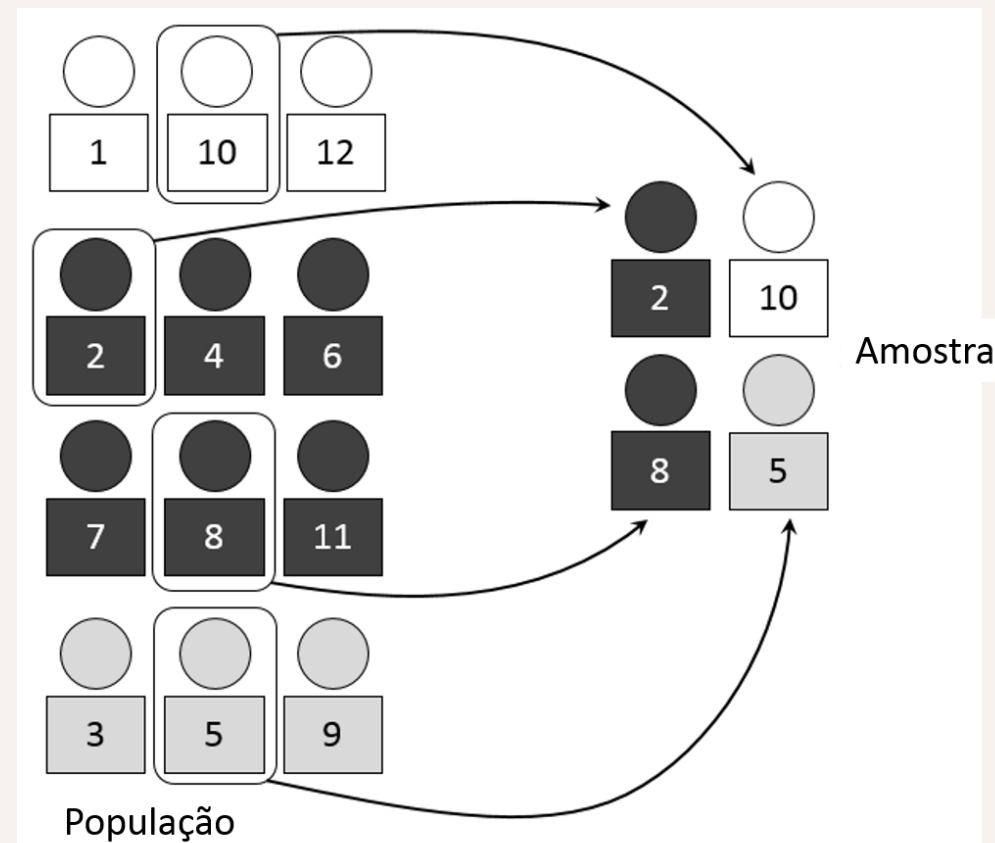
How the Same Voters Voted in the 1932 Election

Electoral Vote	London 1936 Total Vote For State	How the Same Voters Voted in the 1932 Election					Electoral Vote	Roosevelt 1936 Total Vote For State	How the Same Voters Voted in the 1932 Election					Electoral Vote	Landon 1936 Total Vote For State	How the Same Voters Voted in the 1932 Election				
		Rep.	Dem.	Soc.	Other	Did Not Vote			Rep.	Dem.	Soc.	Other	Did Not Vote			Rep.	Dem.	Soc.	Other	Did Not Vote
Ala.	11	3,660	1,218	1,298	5	5	412	126	10,682	371	8,538	50	1							
Ark.	6	2,437	1,431	1,467	18	1	129	117	1,975	286	1,555	33	1							
Cal.	9	2,724	3,338	551	15	1	129	117	1,975	286	1,555	33	1							
Col.	11	1,835	1,500	1,500	15	1	129	117	1,975	286	1,555	33	1							
Conn.	6	15,495	12,214	12,214	15	1	129	117	1,975	286	1,555	33	1							
Del.	3	26,899	22,939	5,176	111	7	12,330	1,140	13,413	2,984	9,115	408	6							
Fla.	7	6,897	1,121	2,011	13	5	344	303	6,420	635	6,724	41	1							
Ga.	12	1,940	1,298	1,617	8	1	1,940	1,298	1,617	8	1,940	1,298	1,617							
Idaho	4	3,653	2,672	698	9	1	683	163	5,611	391	1,989	8	1							
Ill.	26	124,297	81,112	23,885	57	1	4,906	5,125	79,083	14,791	54,412	1,545	6							
Ind.	14	42,895	31,915	7,644	154	49	1,290	1,275	26,663	4,513	20,247	302	22							
Iowa	11	31,871	22,823	6,164	135	1	1,272	1,631	18,644	3,196	13,611	249	7							
Kent.	9	35,488	25,315	6,489	147	15	1,490	1,670	20,254	4,182	14,122	251	11							
Ky.	11	13,385	8,597	7,997	50	14	827	16,392	1,546	13,296	95	6								
La.	10	1,686	1,686	1,686	15	1	789	1,686	1,686	15	1,686	1,686	15							
Leis.	5	11,742	8,619	1,567	35	35	713	781	5,837	635	3,420	41	1							
Mass.	11	1,141	1,141	1,141	15	1	1,141	1,141	1,141	15	1,141	1,141	15							
Mich.	17	87,449	70,567	10,195	330	31	1,213	1,203	28,944	1,141	17,499	744	16							
Min.	16	31,476	31,284	1,665	22	2	2,113	1,863	28,686	1,141	17,402	746	20							
Miss.	11	36,742	22,384	5,934	109	3	972	1,334	29,733	1,699	14,851	511	22							
Mont.	3	1,480	269	394	1	1	4,898	86	5,396	4	798	43	1							
N. H.	13	50,822	31,551	11,491	244	45	2,975	2,098	38,267	4,463	30,608	495	15							
N. J.	14	4,490	1,384	826	16	1	164	1,362	661	2,317	912	100	19							
N. Y.	33	1,083	1,083	1,083	108	7	641	11,770	1,077	9,045	177	2	441							
N. C.	14	9,207	5,04	1,072	21	1	537	2,737	479	1,964	114	114	100							
N. D.	3	1,083	1,083	1,083	27	1	953	1,051	1,051	14	1,051	1,051	14							
Ohio	16	38,672	4,561	8,622	231	17	2,583	1,662	27,632	1,693	10,646	89	1							
Ore.	6	14,686	114,574	35,052	805	45	7,135	4,639	129,277	18,241	99,938	4,101	141							
R. I.	4	4,230	2,297	1,137	15	1	1,988	182	3,664	694	2,679	338	2							
S. C.	8	1,137	1,137	1,137	15	1	1,137	1,137	1,137	15	1,137	1,137	15							
S. D.	3	1,137	1,137	1,137	15	1	1,137	1,137	1,137	15	1,137	1,137	15							
Tenn.	11	14,442	8,931	4,262	29	3	1,630	7,017	15,075	1,289	12,389	53	2							
Tex.	6	11,747	8,931	1,014	29	3	531	19,913	7,661	7,661	7,661	7,661	7,661							
Va.	12	119,084	86,431	20,097	543	115	6,441	5,437	81,114	14,562	56,082	1,340	55							
W. Va.	4	19,481	6,161	1,267	12	1	511	479	13,486	601	479	46	1							
Wash.	8	1,247	219	658	42	14	3,970	7,486	101	5,943	4	4	791							
Wis.	10	4,483	5,472	2,096	42	14	3,970	2,448	371	314	779	122	11							
Wyo.	3	9,851	5,785	2,354	41	31	1,178	5,664	19,829	1,410	13,510	120	3							
Unk.	20	15,841	6,362	1,274	41	31	1,178	5,664	19,829	1,410	13,510	120	3							
Unk.	3	2,841	5,839	851	21	21	139	1,173	5,318	954	3,955	67	8							
Unk.	11	19,222	2,991	2,840	31	18	1,174	410	16,798	1,121	15,156	141	14							



Amostragem estratificada

- Além de ser uma amostragem aleatória, a estratificada busca manter proporções importantes da população
- A amostragem estratificada faz isso ao retirar amostras aleatórias em subconjuntos da população.



Amostragem sistemática

- Todos os elementos da população são conhecidos e uma lógica é usada para selecionar os elementos, por exemplo, selecionar todo décimo elemento.

Amostragem por conveniência

- Usada quando não é possível utilizar a amostragem aleatória, logo, os casos são selecionados pela conveniência do pesquisador.
 - Embora seja o método menos rigoroso, é o mais usado e menos custoso.
 - Tenta-se minimizar os vieses ao buscar assemelhar as características da amostra às da população, quando essas características são conhecidas.
 - **Exemplo:** amostra de estudantes universitários
-

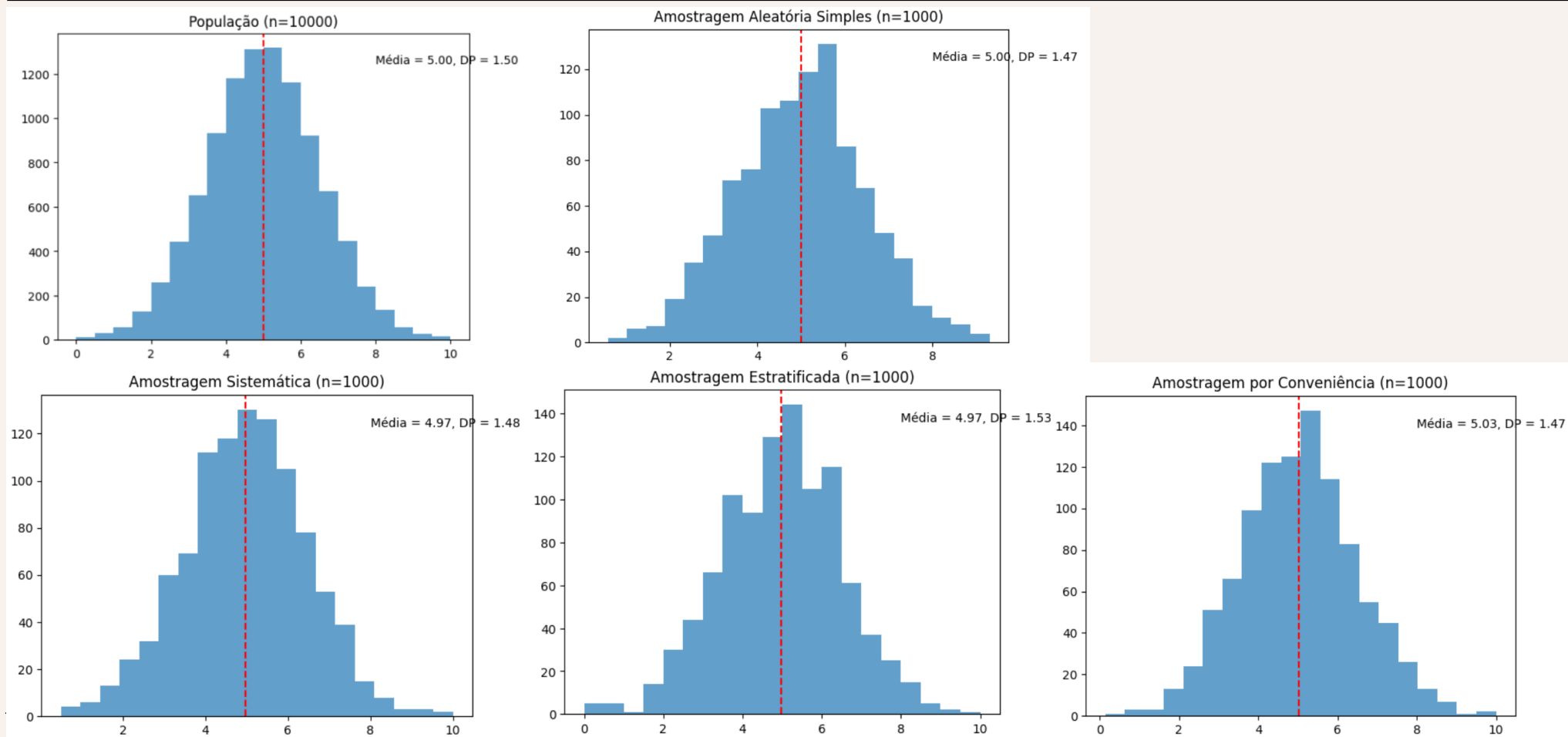
Amostragem por conglomerados

- A população é dividida em grupos chamados de conglomerados e uma amostragem aleatória é realizada entre os grupos.
 - Usada em pesquisas nacionais importantes:
 - PNAD
 - SAEB
 - ...
-

Resumo

Método de Amostragem	Descrição	Exemplo
Amostragem Aleatória Simples	Cada elemento da população tem a mesma probabilidade de ser selecionado.	Sortear números aleatórios para escolher 50 pessoas de um grupo de 500.
Amostragem Sistemática	Os elementos são selecionados a partir de intervalos fixos.	Escolher a cada 10ª pessoa em uma lista ordenada de 1.000 indivíduos.
Amostragem Estratificada	A população é dividida em subgrupos (estratos), e elementos são selecionados de cada grupo.	Selecionar proporcionalmente alunos de diferentes faixas etárias de uma escola.
Amostragem por Conglomerados	A população é dividida em grupos (conglomerados), e alguns grupos são selecionados.	Escolher algumas salas de aula de uma escola e entrevistar todos os alunos dessas salas.
Amostragem por Conveniência	Os elementos são selecionados com base na facilidade de acesso.	Entrevistar as pessoas que estão disponíveis em um shopping.
Amostragem por Quotas	A seleção é baseada em características específicas da população.	Escolher 20 homens e 20 mulheres para uma pesquisa.

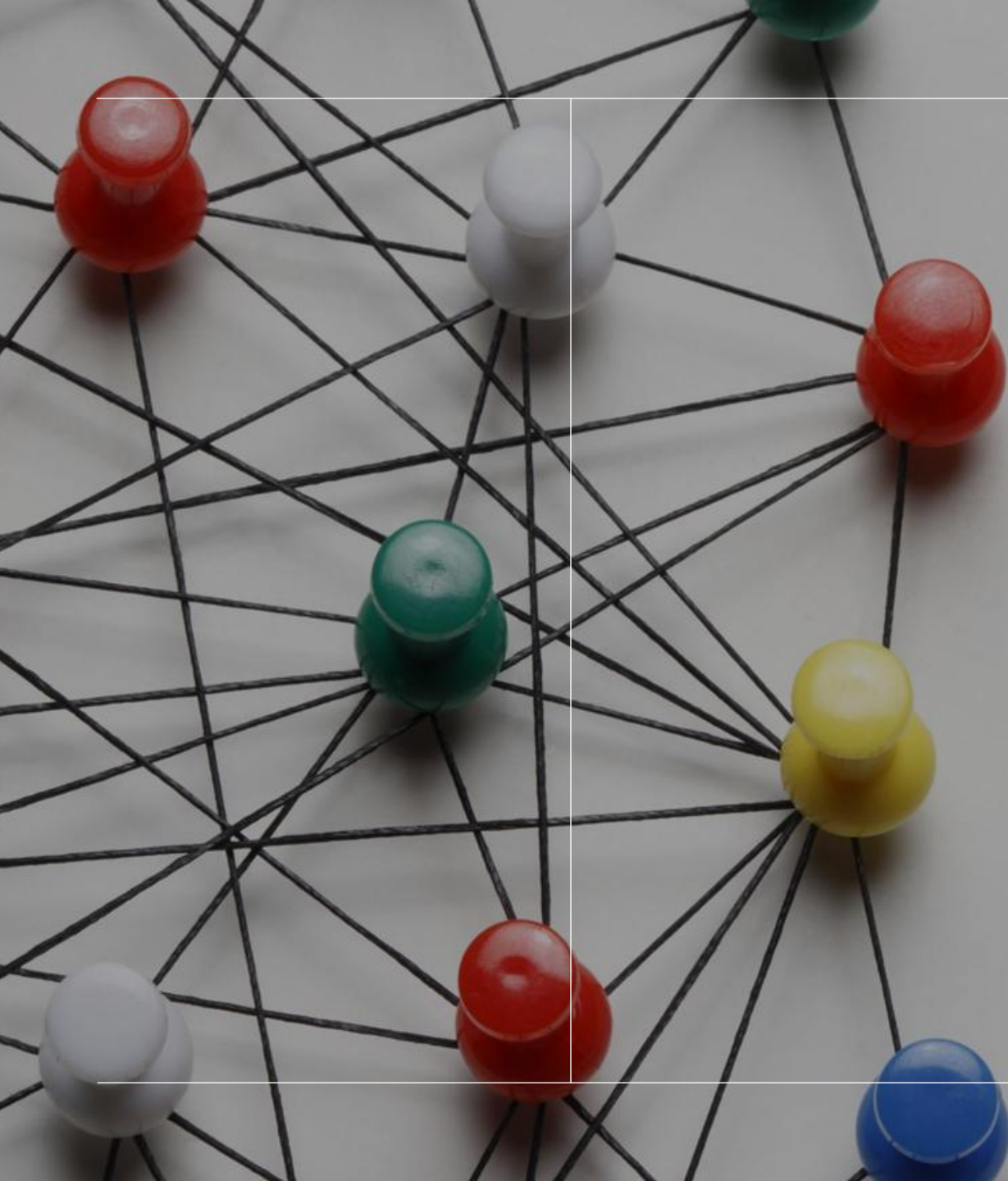
Comparação dos métodos



E aí?

- O que é mais preciso?
- Uma amostra de 200 casos ou 10 amostras de 20 casos?



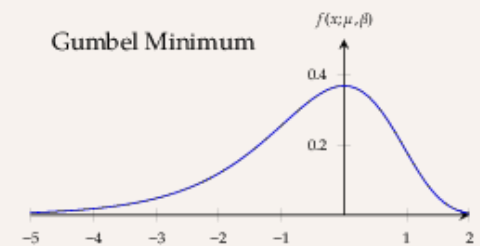
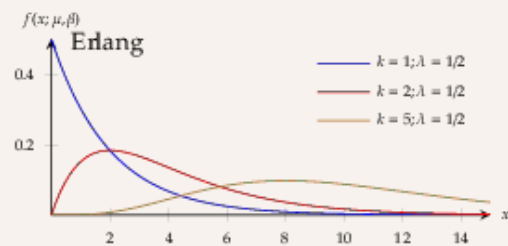
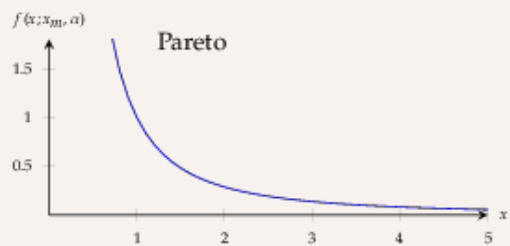
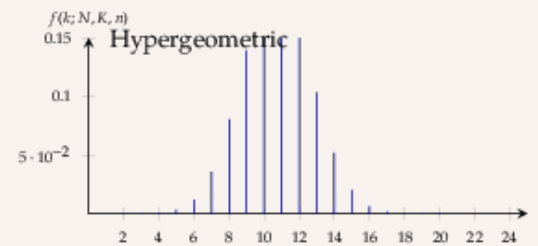
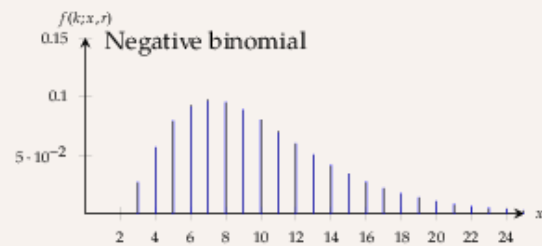
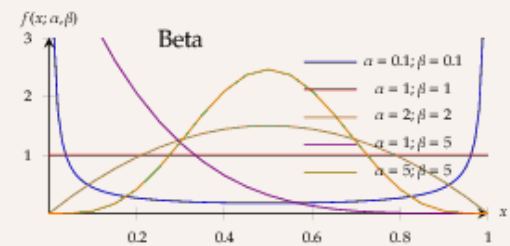
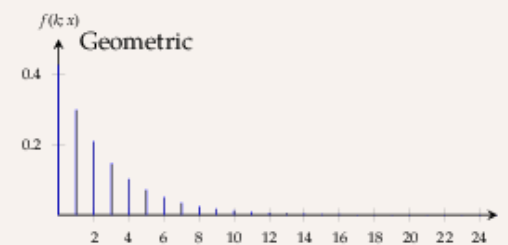
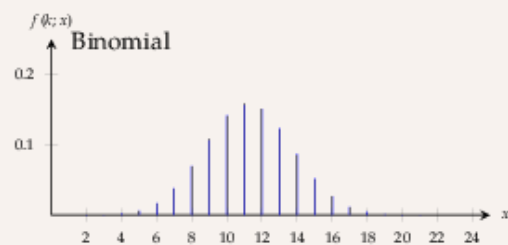
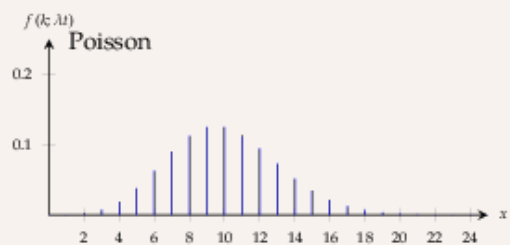


Distribuições
probabilísticas

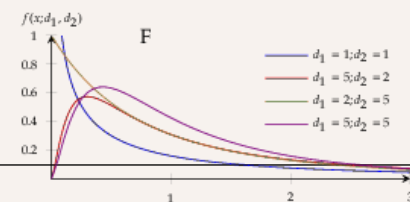
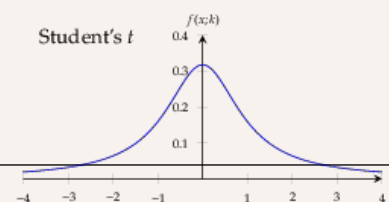
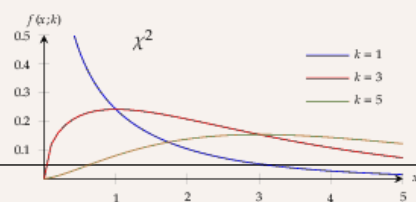
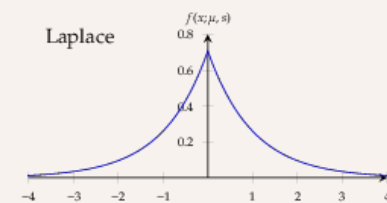
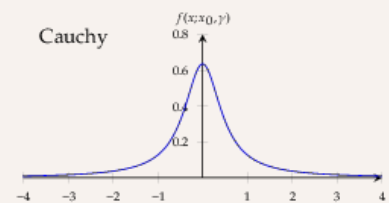
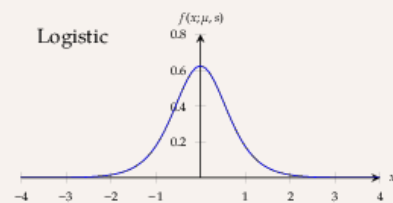
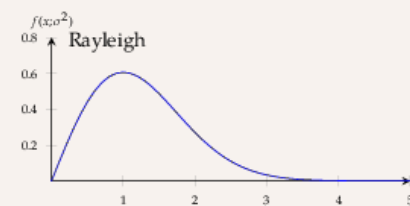
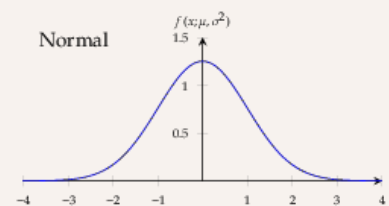
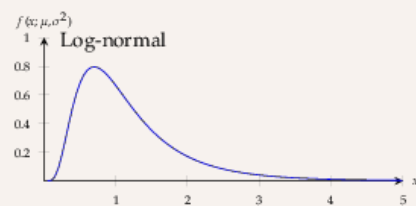
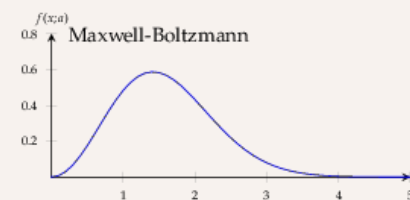
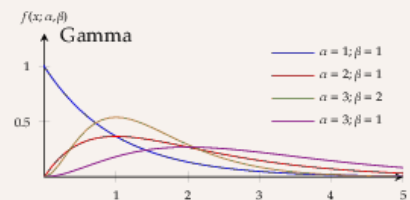
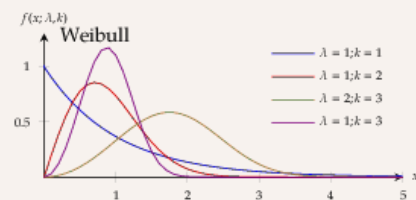
Dois tipos de distribuições

- Existem dois tipos de distribuições:
 - **Teóricas:** usadas para representar padrões de distribuição que, teoricamente, podem existir.
 - **Empíricas:** as distribuições dos dados reais coletados.
 - Podemos visualizar uma distribuição empírica plotando um histograma ou um gráfico de densidade.
 - As distribuições teóricas são úteis para modelar eventos probabilísticos.
 - **Eixo-x** = valores da distribuição
 - **Eixo-y** = probabilidade
 - **Área sob a curva** = 1
-

Várias distribuições



Várias distribuições



Distribuições mais usadas

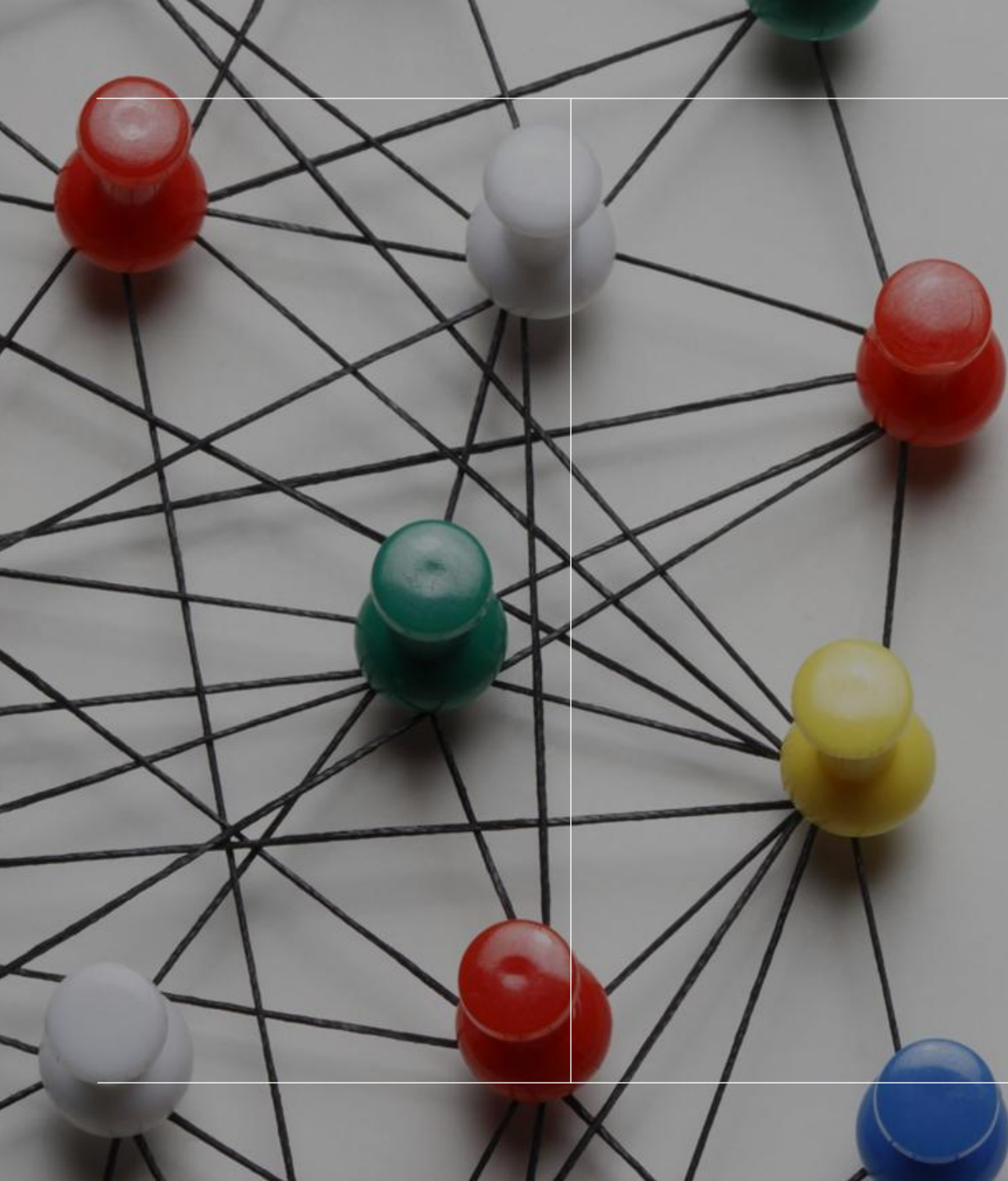
- As distribuições mais usadas na prática são:
 - Normal
 - Exponencial
 - Poisson
 - Logística
 - qui-quadrado
 - t de Student
 - F



Calma!

- Essas distribuições tendem a ser mais importantes para a teoria do que para a prática da estatística e análise de dados. Por isso, não vamos nos aprofundar nelas.





Distribuição normal

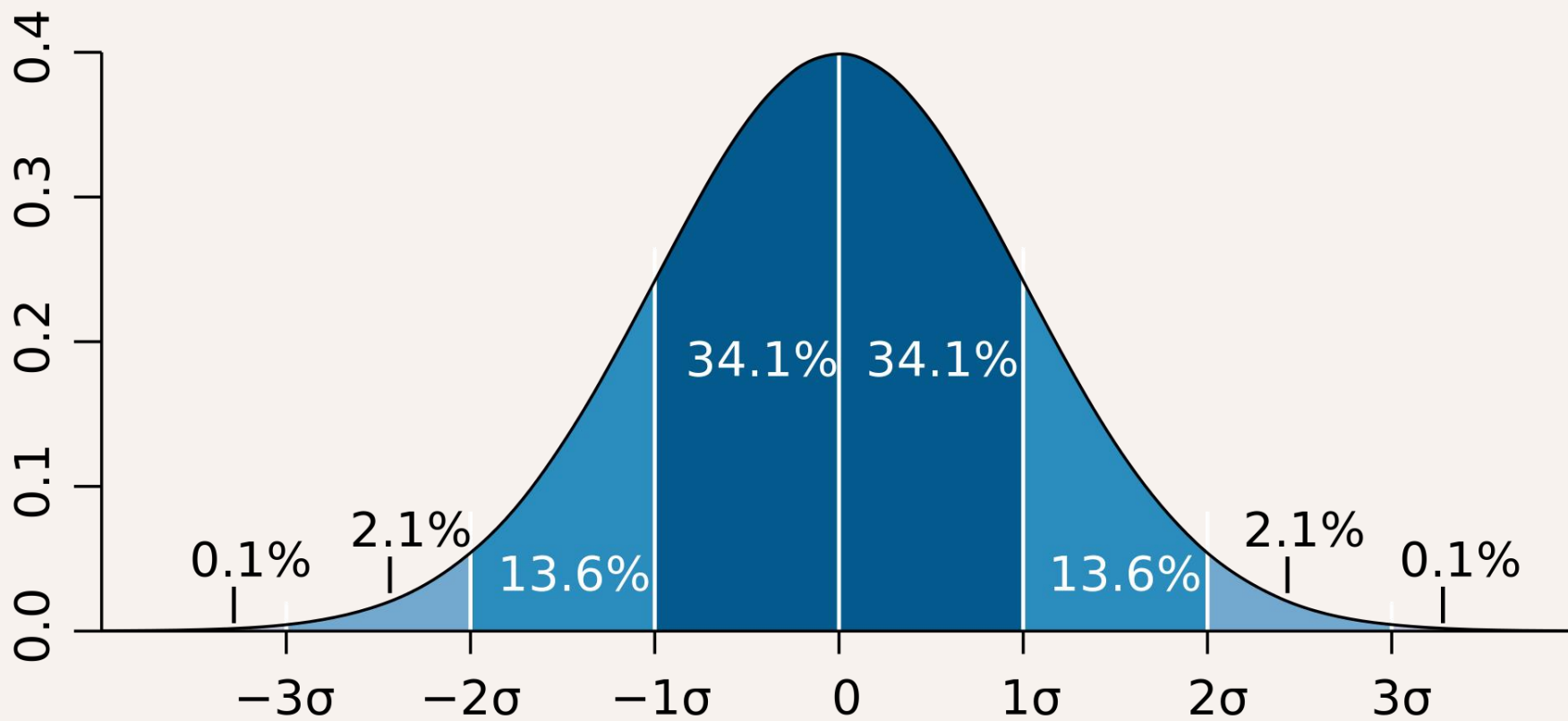
Tábua de Galton



Teórica e empírica

- O nome distribuição normal (Gaussiana) pode se referir a duas coisas:
 - A distribuição de probabilidades teórica chamada distribuição normal padrão (DNP);
 - Dizemos que uma variável no nosso conjunto de dados tem distribuição normal quando sua forma é semelhante à DNP (uma distribuição normal empírica).
-

DNP



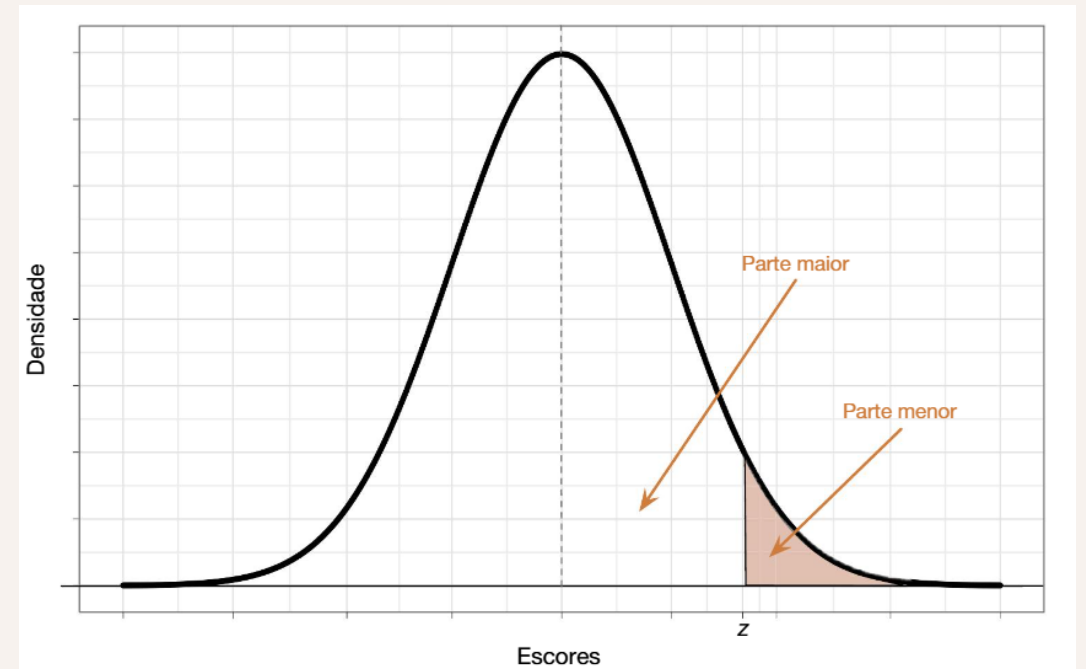
Teórica e empírica

- Para cada valor da DNP há uma probabilidade correspondente, i.e.

$$x_i \rightarrow y_i = p(x_i)$$

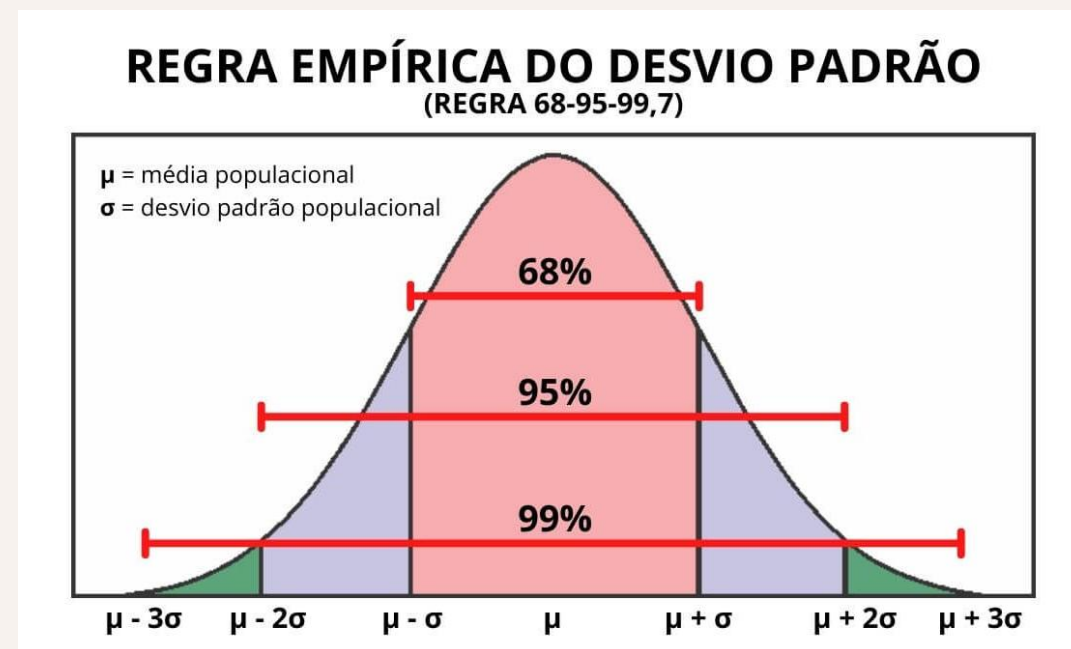
Características

- Forma de sino
- Centrada na média e simétrica em torno dela
- **Média = 0 e desvio-padrão = 1**
- Cada escore no eixo-x é chamado de escore Z, uma medida de desvios-padrão da média.
- Para cada escore Z existe uma região abaixo e acima, que é a probabilidade de um escore Z menor que ou maior que aquele.



Regra empírica

- Uma característica da DNP é a regra empírica, que especifica proporções da curva para intervalos de desvio-padrão acima e abaixo da média:
 - ± 1 DP da média = 68% da curva
 - ± 2 DP da média = 95% da curva
 - ± 3 DP da média = 99% da curva
- Na realidade, qualquer intervalo pode ser computado usando os escores z. As probabilidades acima e abaixo da curva para cada escore z já estão tabeladas.



Escores Z

- Para usar as propriedades da distribuição normal com os valores de uma variável, posso transformar os escores da minha variável em escores z:

$$z_i = \frac{x_i - \bar{x}}{DP}$$

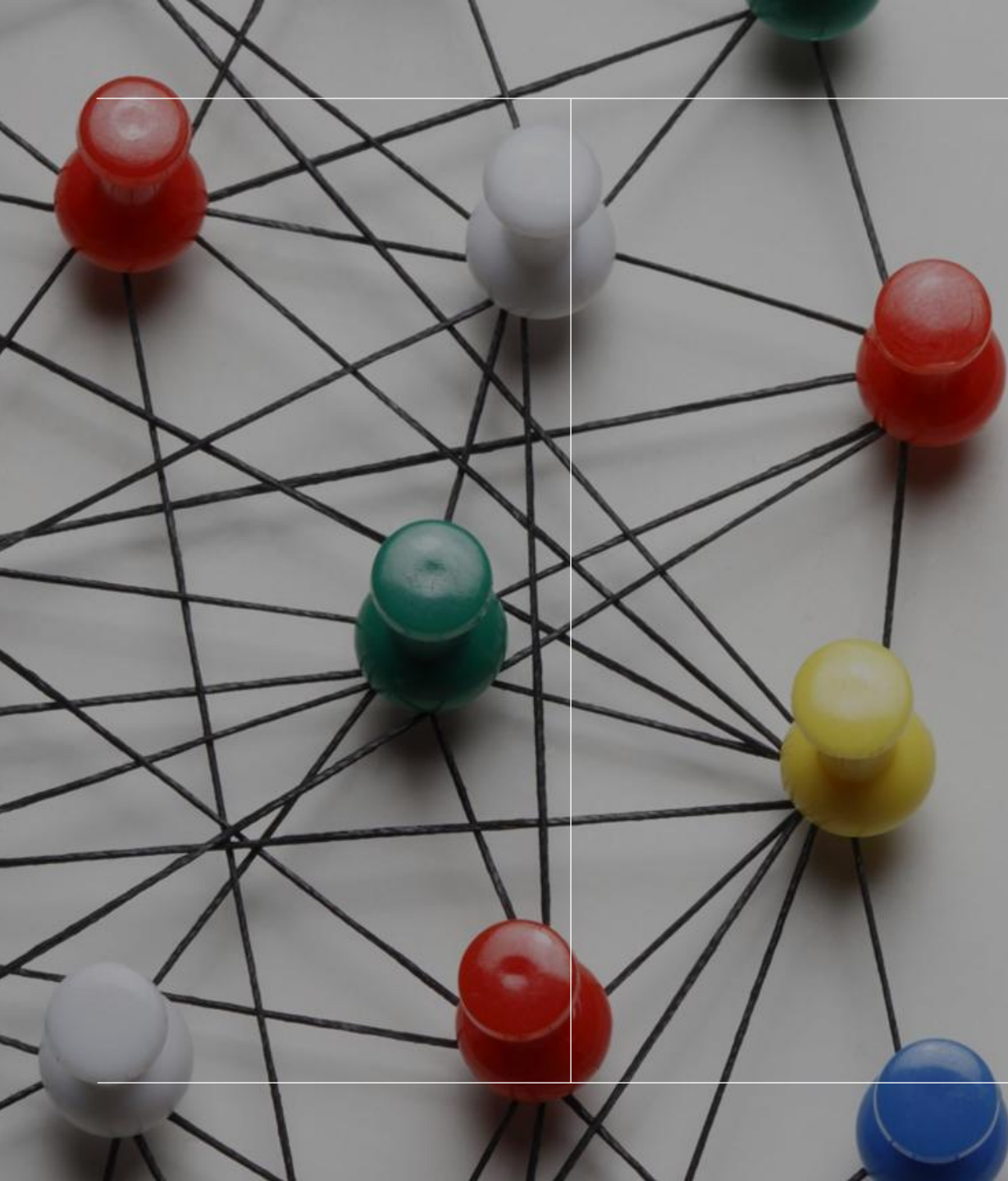
- Assim, os escores z representam o afastamento da média na DNP, cuja média é 0 e desvio-padrão é 1.

	Variável	Desvios	Escore Z
	1	-3,86	-1,33
	4	-0,86	-0,30
	7	2,14	0,74
	8	3,14	1,08
	2	-2,86	-0,99
	3	-1,86	-0,64
	9	4,14	1,43
Média	4,86		0
DP	2,90		1

Exercício

- Use a tabela de escores Z para encontrar as probabilidades dos escores Z correspondentes!

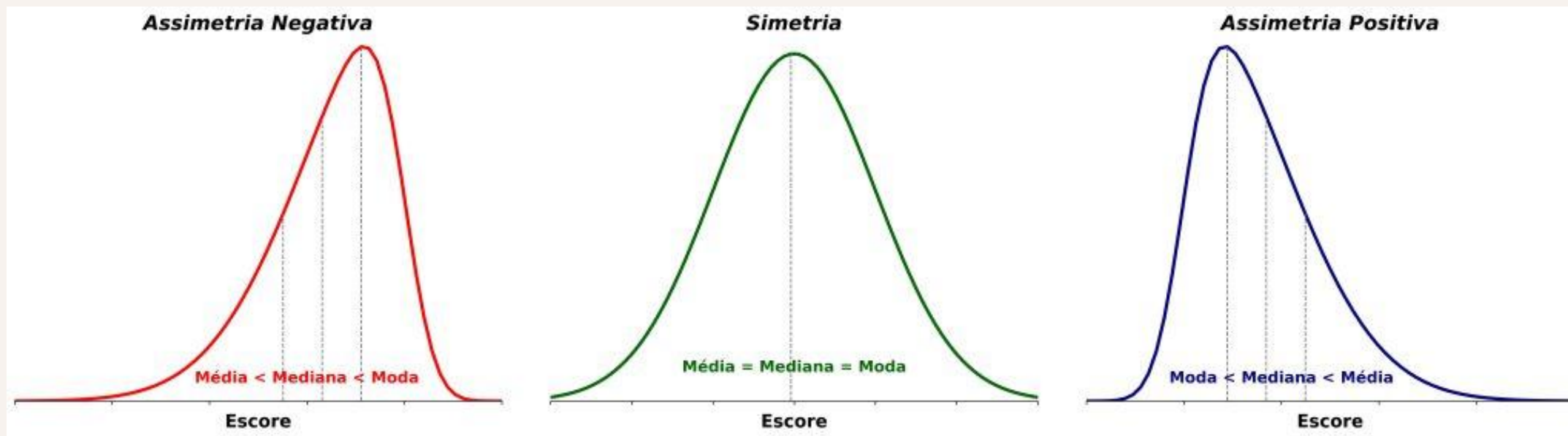
	Variável	Desvios	Escore Z
	1	-3,86	-1,33
	4	-0,86	-0,30
	7	2,14	0,74
	8	3,14	1,08
	2	-2,86	-0,99
	3	-1,86	-0,64
	9	4,14	1,43
Média	4,86		0
DP	2,90		1



Desvios da normalidade

Assimetria

- A distribuição normal é simétrica em torno da média, *i.e.*, tem a mesma proporção de dados acima e abaixo da média (50-50%). Uma distribuição assimétrica viola essa característica.

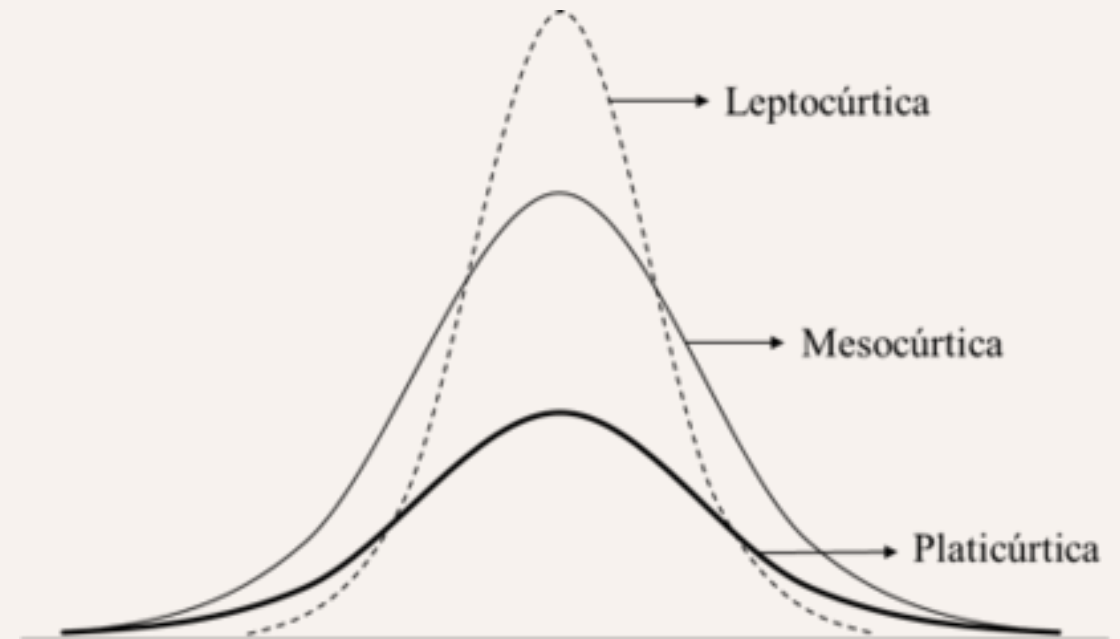


Métrica

- A assimetria (desvio da simetria, *skewness*) pode ser calculada e assumir valores positivos, negativos ou zero:
 - **Assimetria = 0**: distribuição simétrica;
 - **Assimetria > 0**: assimetria positiva;
 - **Assimetria < 0**: assimetria negativa.
-

Curtose

- A distribuição normal é mesocúrtica, i.e.:
 - Nem pontuda demais: dados muito concentrados na média, desvio-padrão muito pequeno;
 - Nem larga demais: dados muito espalhados da média, desvio-padrão grande.



Métrica

- A curtose (*kurtosis*) pode ser calculada e assumirá valores positivos, negativos ou zero:
 - **Curtose = 3**: distribuição mesocúrtica;
 - **Curtose > 3**: distribuição leptocúrtica;
 - **Curtose < 3**: distribuição platicúrtica.
 - Em algumas implementações, calcula-se o excesso de curtose:
 - **Curtose = 0**: distribuição mesocúrtica;
 - **Curtose > 0**: distribuição leptocúrtica;
 - **Curtose < 0**: distribuição platicúrtica.
-

Como saber se a minha distribuição é normal?

- Podemos:
 - Avaliar o histograma
 - Verificar os valores de assimetria e curtose
 - Plotar um gráfico-QQ
 - Realizar o teste de normalidade

Gráfico-QQ

- Gráfico Quantil-Quantil:
 - Converte os valores da minha variável para percentis (eixo-x)
 - Converte os escores Z da DNP para percentis (eixo-y)
 - Plota os eixos em um gráfico
 - Se os pontos recaem sob a linha diagonal, a distribuição é normal

 - Exemplo no JASP!
-

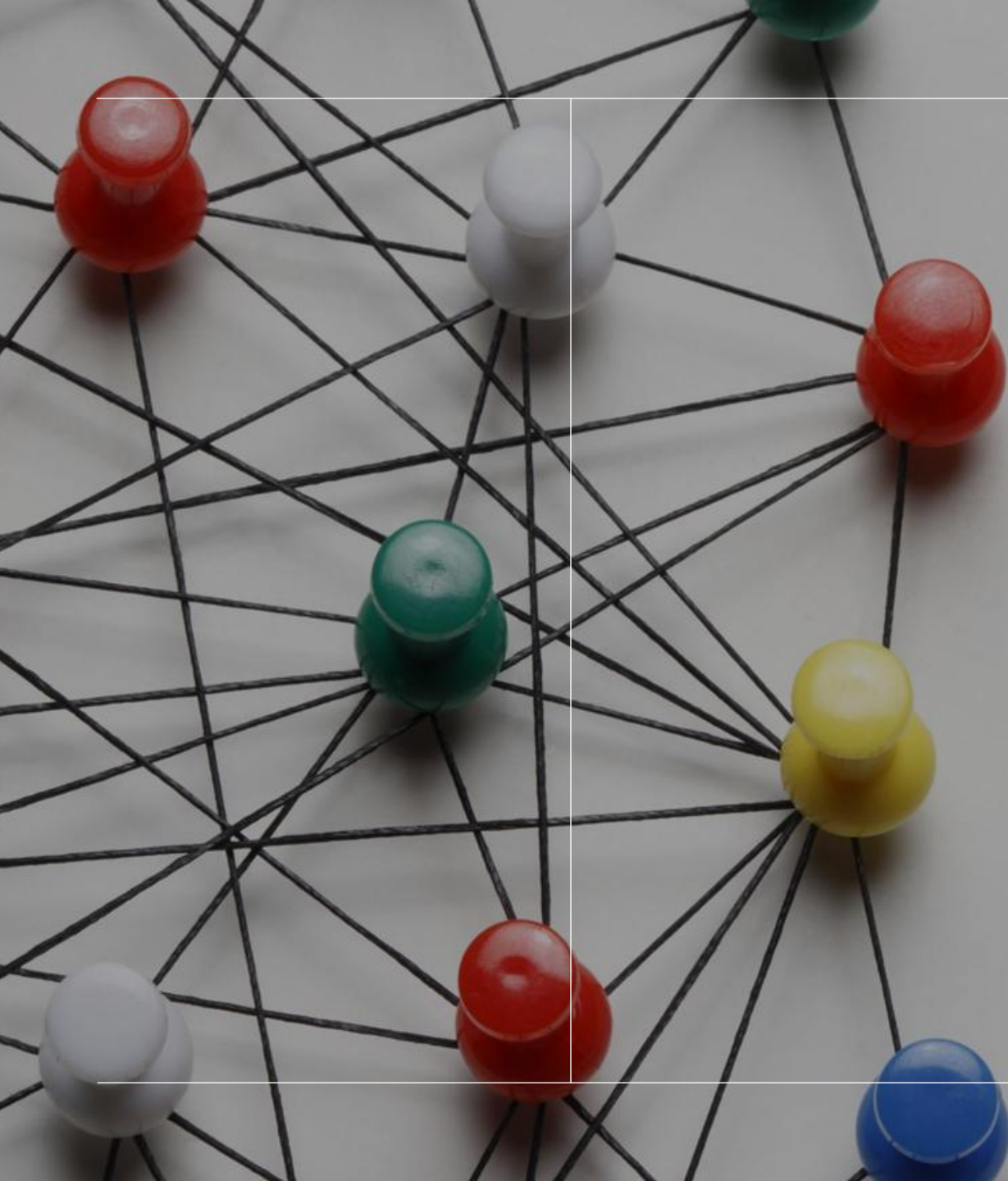
Testes de normalidade

- Avaliam quão diferente a minha distribuição é de uma distribuição normal:
 - Teste de Kolmogorov-Smirnov ou Shapiro-Wilk
 - São muito usados, mas não deveriam
 - Se $p < 0,05$, então a distribuição não é normal
 - Esses testes tem muitos problemas... amostras grandes sempre dão $p < 0,05$

 - Exemplo no JASP!
-

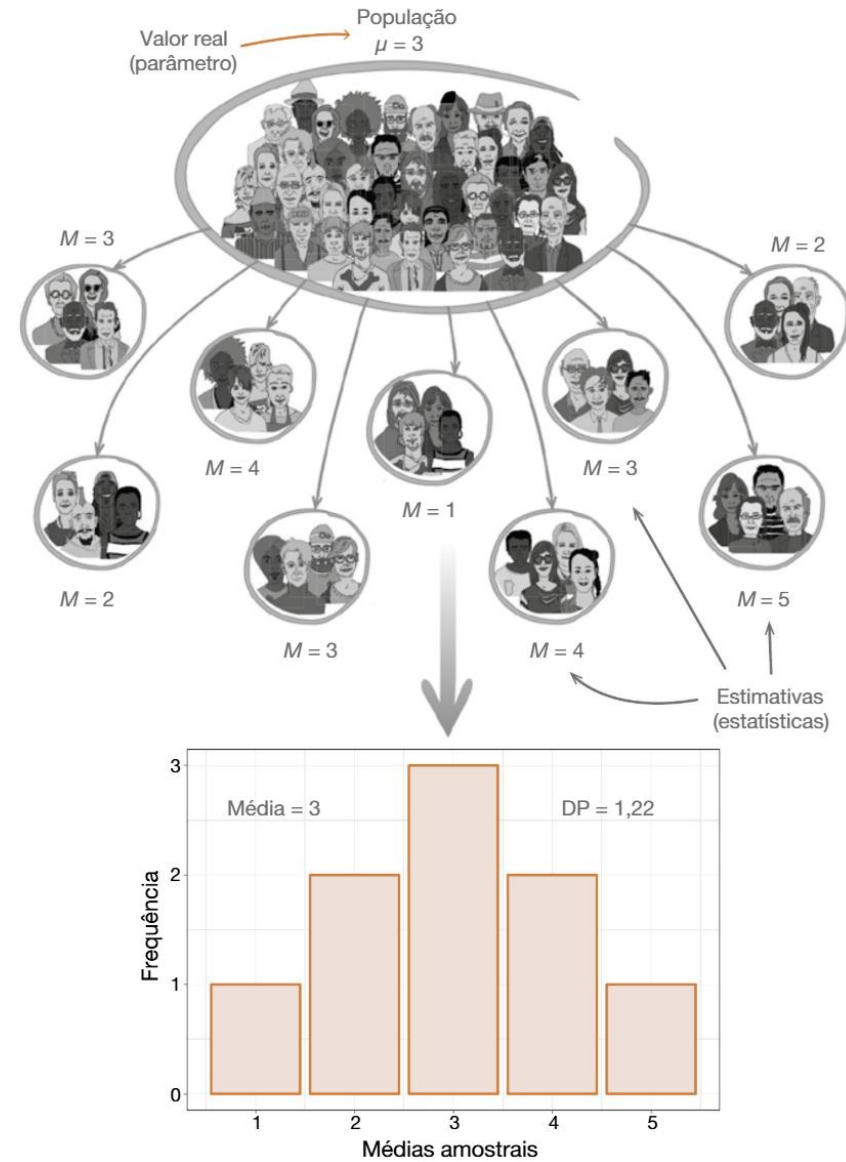
É importante saber se a minha dist. é normal?

- Uma vez era, hoje nem tanto!
 - É importante saber se existem grandes desvios da normalidade, mas isso não vai nos impedir de calcular estatísticas.
 - Não é a normalidade dos dados que importa, mas da distribuição amostral (em breve).
 - Métodos computacionais como *bootstrap* podem contornar a não-normalidade;
-



Distribuição amostral

Experimento mental



Experimento mental

- Vamos fazer um experimento mental:
 1. Retiramos uma amostra aleatória de **30 casos** ($n = 30$) de uma certa população
 2. Calculamos uma estimativa nessa amostra (e.g., média)
 3. Repetimos os passos 1 e 2 **muitas vezes** (k)
 4. Montamos um histograma com as estimativas
 5. A distribuição representada pelo histograma com as estimativas é a distribuição amostral
-

Exemplo

- Vamos ver isso!
-

Normalidade

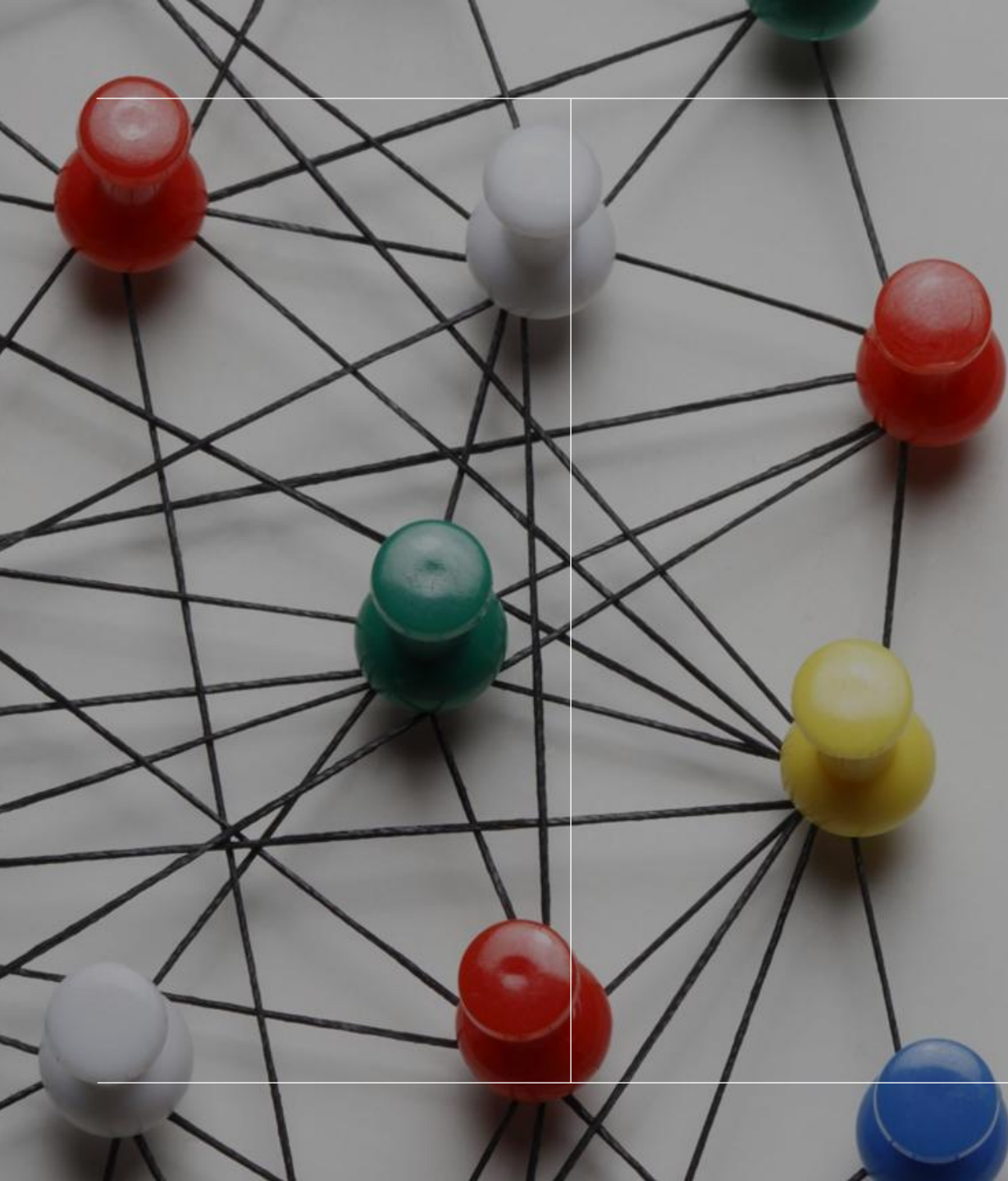
- Na medida em que as amostras retiradas se tornam maiores, a distribuição amostral da estatística calculada nessa amostra tende a ter uma forma de sino (normal);
 - Isso acontece mesmo se tirarmos as amostras de uma distribuição que não seja normal;
-

Erro

- A média das estatísticas usadas para a distribuição amostral é a melhor estimativa possível do parâmetro populacional.
 - O desvio-padrão da distribuição amostral indica a variabilidade entre amostras na estimativa do parâmetro, i.e., o grau médio de erro, chamado de erro-padrão.
-

É pura Teoria!

- A distribuição amostral é uma ideia, uma teoria. Nunca vamos montar uma de fato (exceto no processo de *bootstrap*).
 - Usamos essa ideia para entender que nossa estimativa é uma dentre muitas possíveis (podemos tirar diferentes amostras da mesma população) e, com amostras maiores, sabemos que estamos mais perto do valor verdadeiro do parâmetro.
-



Teorema do Limite Central

Fundamental!

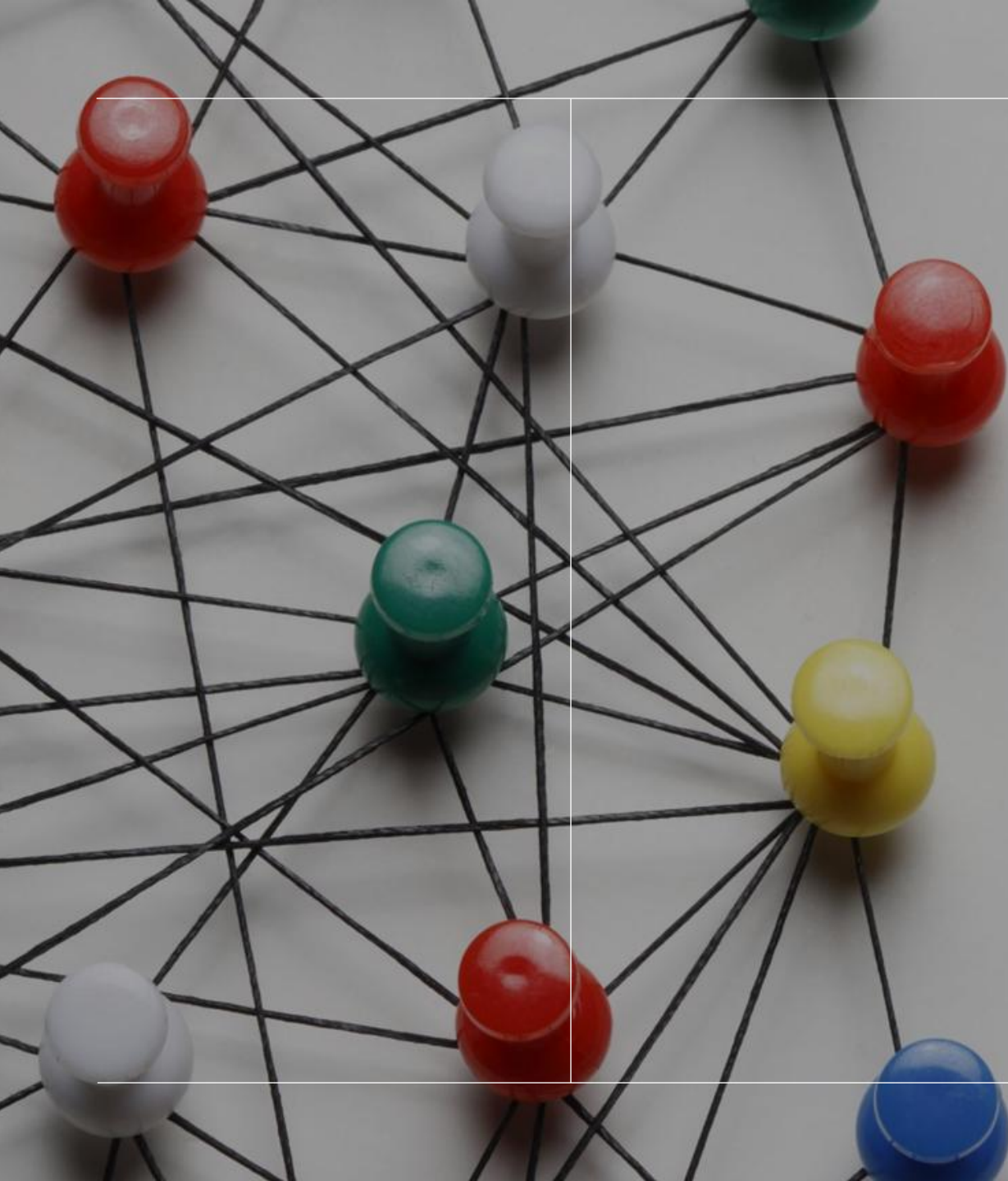
- O teorema do limite central (TLC) é um dos resultados matemáticos mais importantes na estatística.
 - Ele diz que se retirarmos infinitas amostras aleatórias de uma população, cada um de tamanho $n \geq 30$, e calcularmos uma estatística em cada amostra, reunindo essas estatísticas em uma distribuição amostral, então teremos uma distribuição amostral com forma normal, cuja média será uma boa estimativa do parâmetro na população.
-

Por partes ...

- Retirar infinitas amostras aleatórias de uma população ($n \geq 30$) +
 - Calcular uma estatística em cada amostra (e.g., uma média) +
 - Reunir essas estatísticas em uma distribuição +
-
- = distribuição amostral com forma normal
 - = média da distribuição amostral é igual ao parâmetro na população
-

Implicação

- Se tivermos uma amostra grande o suficiente, não vamos precisar nos preocupar com a normalidade da distribuição amostral!



Erro-padrão

Definição

- O erro-padrão é o desvio-padrão da distribuição amostral.
 - Também é chamado de erro amostral, *i.e.*, a diferença média entre diferentes amostras na estatística calculada. Em outras palavras, a métrica que calculamos varia de amostra para amostra.
 - Se fizéssemos o processo de formar a distribuição amostral e calculássemos o desvio-padrão dessa distribuição, teríamos o erro-padrão.
-

Cálculo

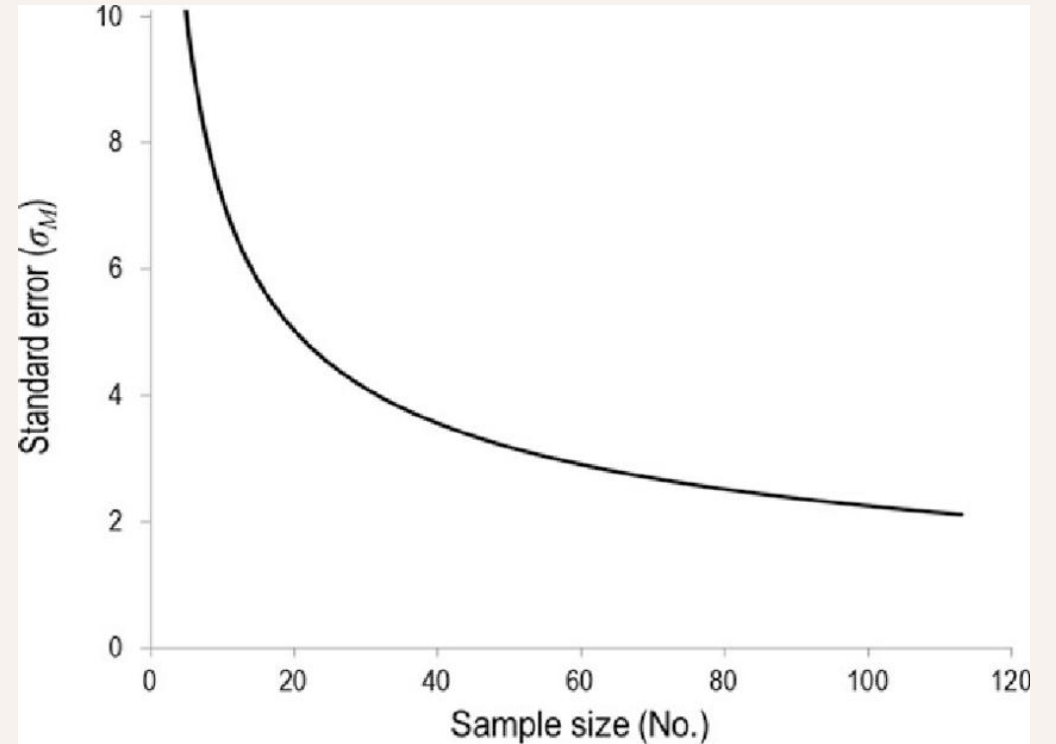
- Para evitar o processo de formar uma distribuição amostral (que seria muito custoso), existe um cálculo direto do erro-padrão, onde:
- EP = Erro-Padrão
- DP = Desvio-padrão da amostra
- N = tamanho da amostra

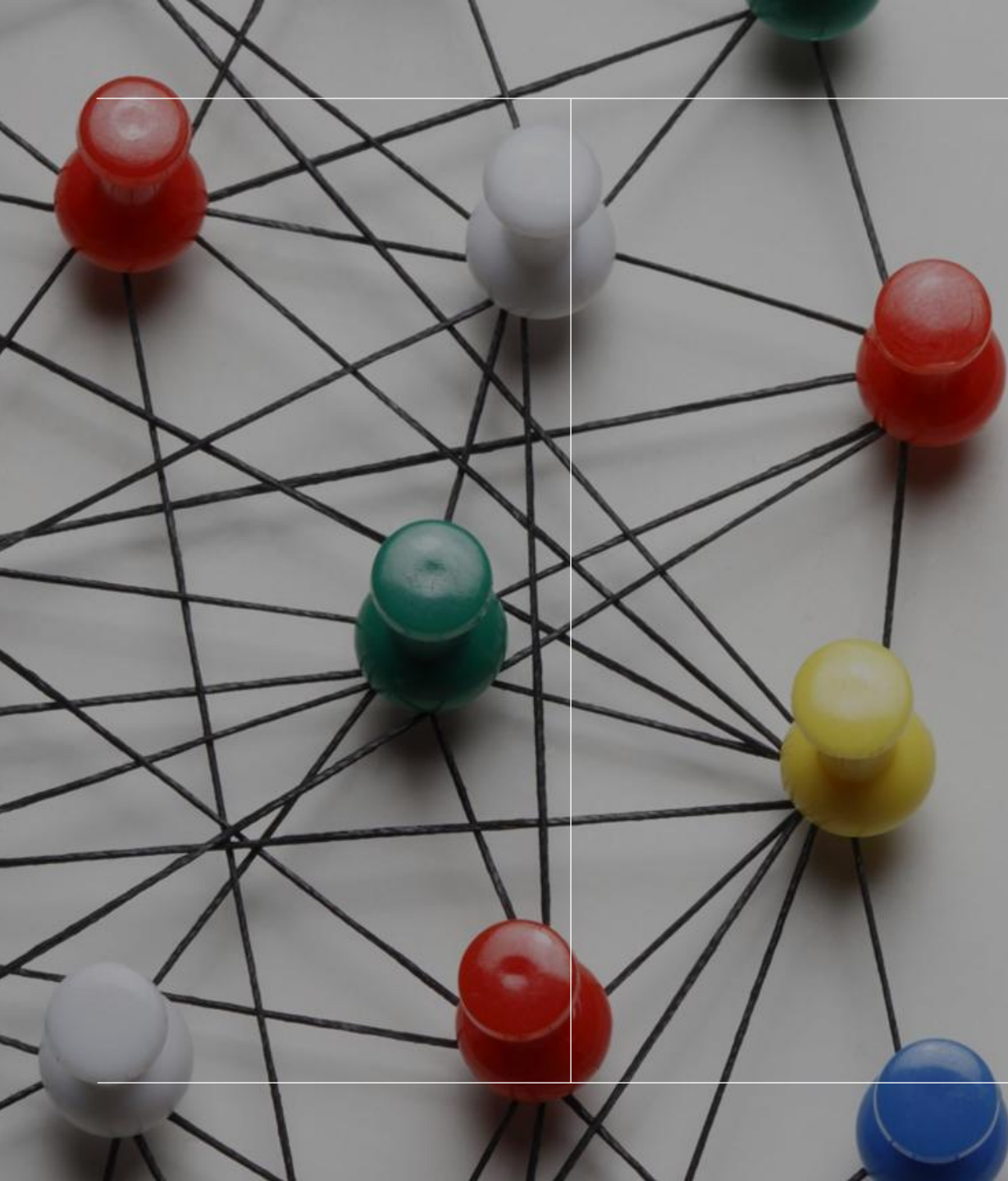
$$EP = \frac{DP}{\sqrt{N}}$$

Característica fundamental do EP

- O erro-padrão caí na medida em que o tamanho da amostra aumenta.
- Tamanhos de amostra grande produzirão erros-padrão muito pequenos.

$$EP = \frac{DP}{\sqrt{N}}$$





Intervalo de confiança

Estimativa pontual e intervalar

- A estatística calculado em nossa amostra é uma estimativa pontual do parâmetro da população ... mas, e se tivéssemos como estimativa um intervalo de valores dentro do qual o parâmetro pode estar?
 - Essa estimativa intervalar chama-se intervalo de confiança!
-

Exemplo:

- Vamos imaginar o seguinte caso:
 - Temos uma amostra de 498 estudantes brasileiros que se inscreveram para cursar uma pós-graduação em uma universidade americana, cada um com um escore do TOEFL, com escores que podem variar de 0 a 120.
 - Queremos saber o escore médio do TOEFL na população de estudantes brasileiros que se inscrevem para uma universidade americana
 - Calculamos a média na amostra: $M = 107,20$
-

Estatísticas

- Calculamos a média na amostra: $M = 107,20$
 - Calculamos o desvio-padrão: $DP = 6,09$
 - Calculamos o erro-padrão: $\frac{6,09}{\sqrt{498}} = 0,273$
-

Sabemos ...

- Que entre $-1,96$ e $1,96$ escores-Z da distribuição normal padrão temos 95% da área sob a curva, ou seja, temos uma probabilidade de 95% de encontrar um escore entre $-1,96$ e $1,96$ se retirarmos um valor qualquer de uma distribuição normal padrão.
 - Que a distribuição amostral é uma distribuição normal.
 - Logo, podemos usar o erro-padrão para calcular um intervalo de 95% na distribuição amostral em torno da estatística calculada na amostra.
-

Cálculo do intervalo para média

- Calculamos esse intervalo da seguinte forma:

$$Limite_{superior} = \bar{x} + 1,96 \times EP$$

$$Limite_{inferior} = \bar{x} - 1,96 \times EP$$

- No exemplo:
 - Média: 107,20
 - LS: 107,73
 - LI: 106,66
-

Interpretação frequentista

- Um intervalo de confiança de 95% significa que em 95 de 100 vezes que retiramos uma amostra, calcularmos uma estatística e calcularmos um intervalo, o verdadeiro valor do parâmetro estará dentro do intervalo. Em outras palavras, somente em 5 vezes o valor do parâmetro não estará dentro do intervalo.
 - **Interpretações erradas:**
 - O intervalo de confiança indica uma probabilidade de 95% do intervalo conter o valor do parâmetro
 - O intervalo de confiança indica que só há 5% de chance de erro aleatório
-

Interpretação frequentista

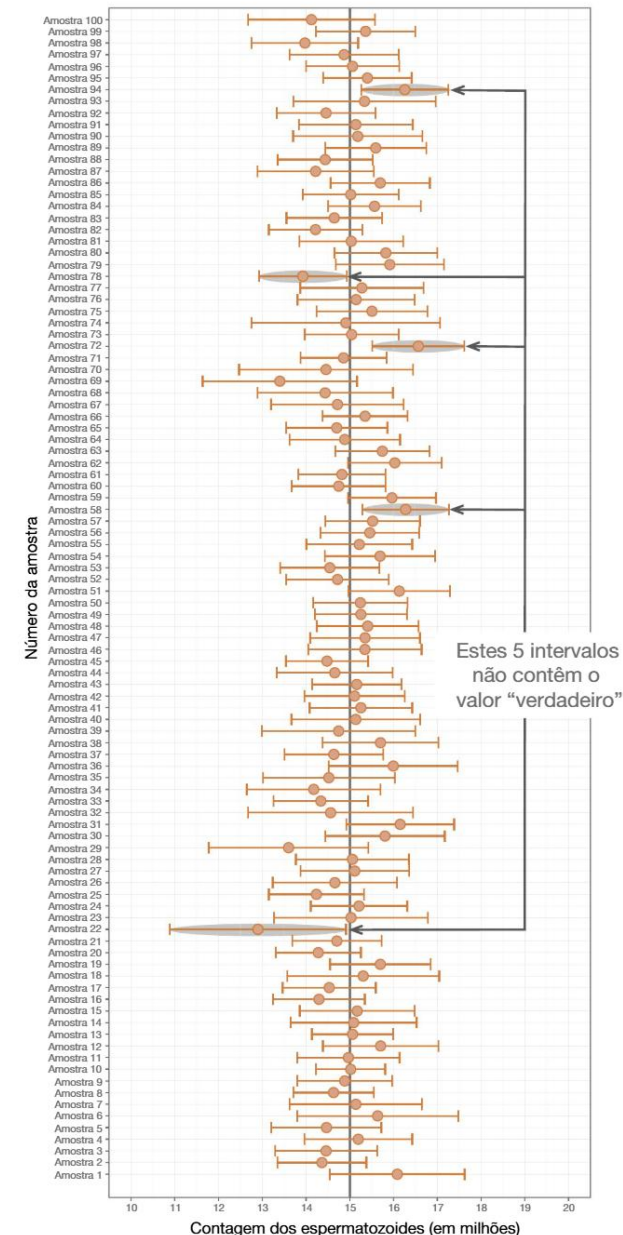


Figura 2.9 Intervalos de confiança das contagens de espermatozoides da codorna japonesa (eixo horizontal) para 100 amostras diferentes (eixo vertical).

Características

- Podemos calcular intervalos de confiança para qualquer estatística. Versão generalizada da fórmula: $IC = b \pm z \times EP$
 - Os intervalos são afetados pelo valor de Z. Quanto maior, mais amplo fica o intervalo:
 - 90%: $z = 1,68$
 - 95%: $z = 1,98$
 - 99%: $z = 2,57$
 - 99,9%: $z = 3,29$
 - Os intervalos são afetados pelo erro padrão. Quanto mais erro, mais amplo é o intervalo.
-

Características

- Quanto mais estreito o IC, mais precisa é nossa estatística.
 - Intervalos que cruzam o zero sugerem um grau de erro muito elevado na estimativa, exceto em alguns poucos casos.
-

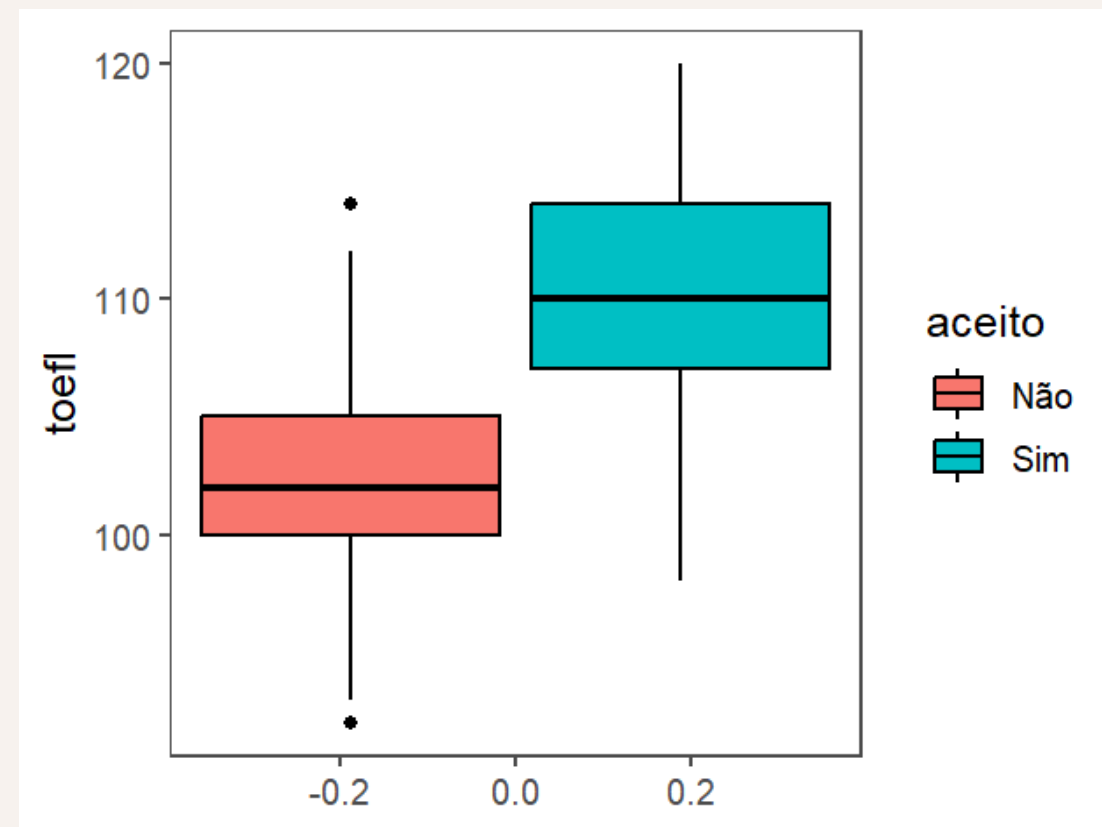
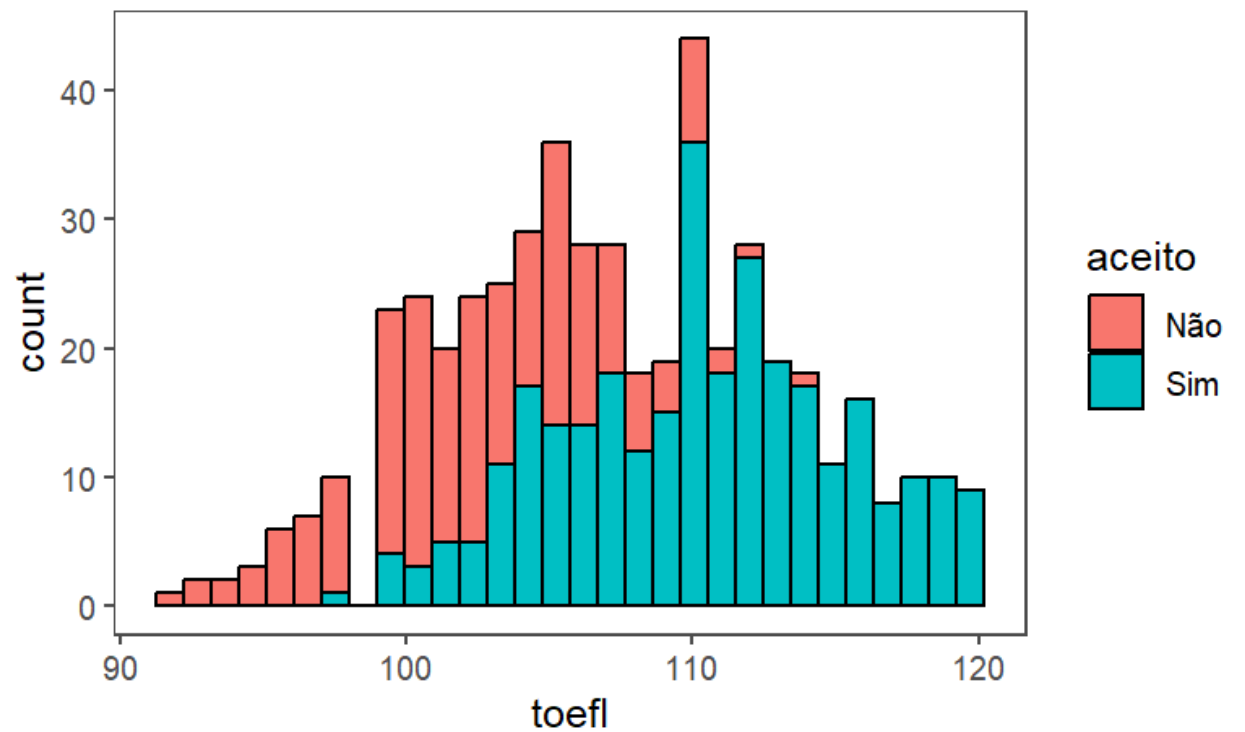
Exemplo de uso dos intervalos

- Existe diferença nos escores do TOEFL entre um grupo de estudantes que foi aprovado para uma universidade americana e um grupo que não foi?
 - Média Aprovados = 110,32
 - Média Não-Aprovados = 102,46
 - Existe uma diferença matemática, mas essa é uma diferença estatística? As distribuições dos dois grupos são mais diferentes do que semelhantes entre si?
-

Diferenças

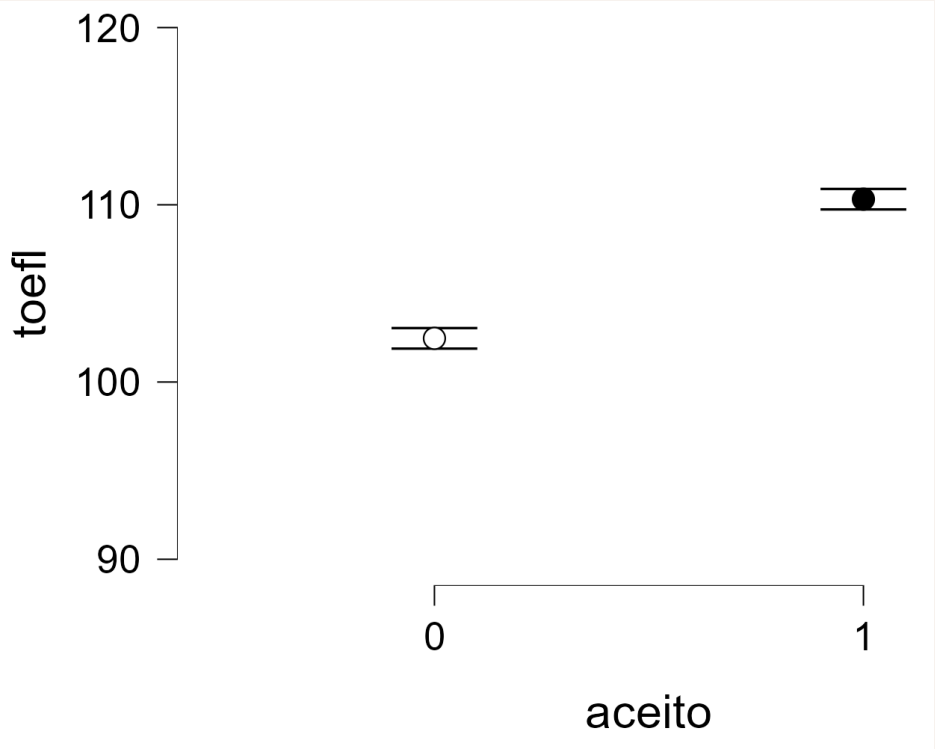
- **Diferença matemática:** dois números, quando subtraídos, tem resultado diferente de 0, i.e., $x - y \neq 0$
 - **Diferença estatística:** uma diferença entre dois números que é suficientemente grande ao ponto de não poder ter ocorrido ao acaso, i.e., os números são métricas calculadas de distribuições distintas, vindas de populações distintas. Muitas diferenças podem ocorrer ao acaso quando lidamos com amostras.
-

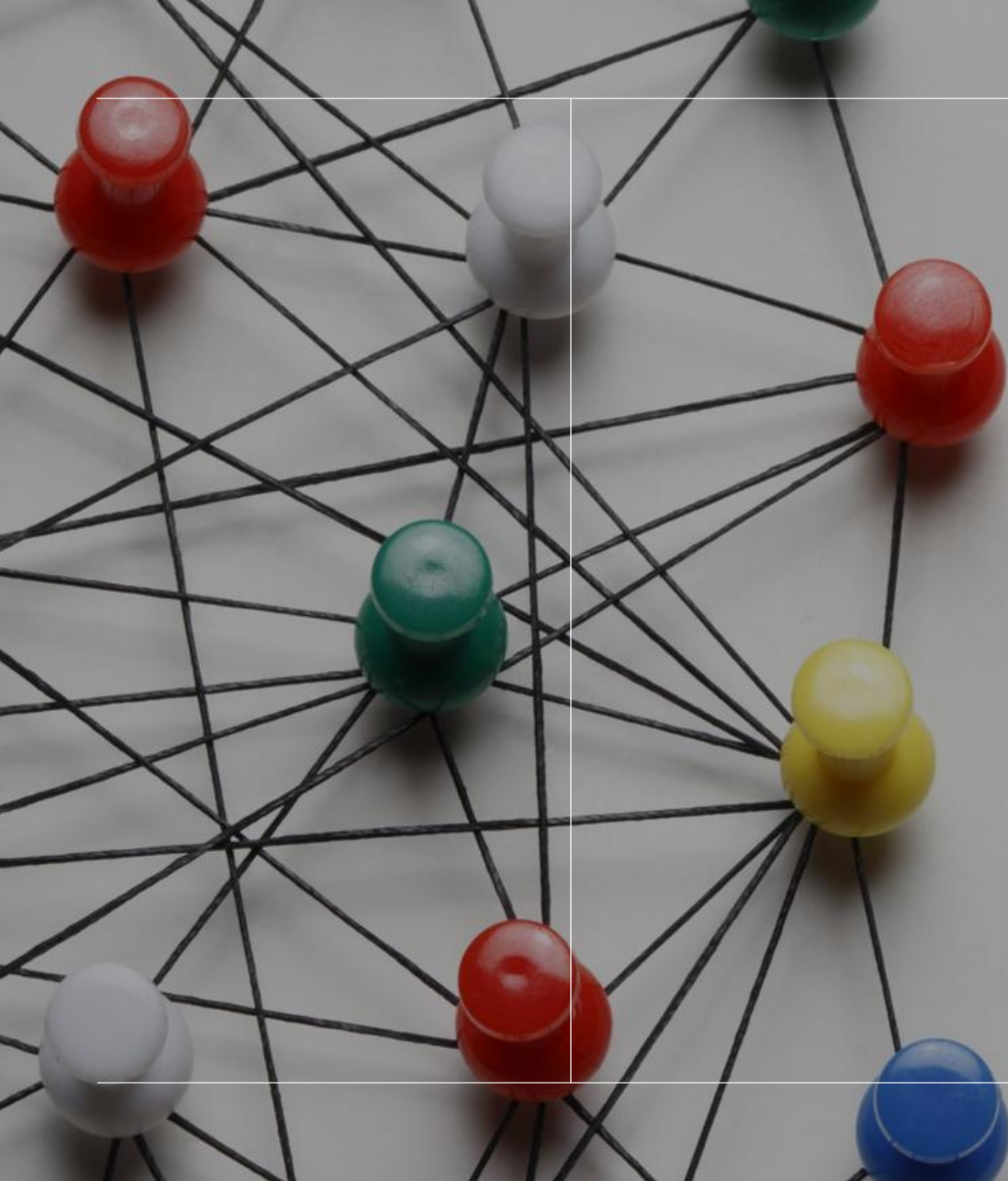
Comparações



Intervalos

Grupo	Média	IC 95% [LI; LS]
Aceito (<i>n</i> = 300)	102,47	101,89; 103,04
Não-aceito (<i>n</i> = 198)	110,32	109,74; 110,89





Testagem de hipóteses

Propósito

- O Teste de Significância da Hipótese Nula (TSHN ou NHST) é um processo desenvolvido para testar hipóteses verificando se nossas estatísticas representam algo real ou se devem ao acaso (variação amostral).
 - As estatísticas operacionalizam nossas hipóteses, i.e., trazem elas do mundo teórico para a realidade física/empírica.
-

Hipóteses

- No TSHN sempre vão existir dois tipos de hipóteses:
 - Hipótese nula **h_0** : representa a nulidade, ausência de efeito, acaso.
 - Exemplo: tomar ômega 3 não está associado ao desempenho cognitivo.
 - Hipótese alternativa **h_1** : representa um efeito.
 - Exemplo: tomar ômega 3 está associado ao desempenho cognitivo.
-

Hipóteses

- As hipóteses podem ser:
 - **Unilaterais/unicaudais:** especificam uma direção para o efeito, por exemplo, tomar ômega 3 aumenta o desempenho cognitivo;
 - **Bilaterais/bicaudais:** não especificam uma direção para o efeito, por exemplo, tomar ômega 3 está associado ao desempenho cognitivo.
-

Tipos de hipóteses

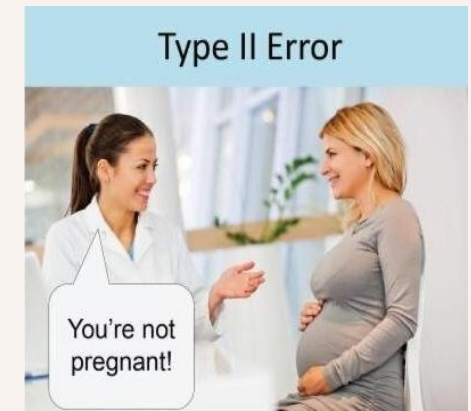
*Os testes de hipóteses, em geral, se referem às relações entre as variáveis. Cada "estatística" na tabela é um só valor que representa a relação entre duas variáveis.

A questão é: esse valor pode ter acontecido ao acaso porque nossa amostra era diferente da população ou ele indica algo real, uma padrão que acontece na população?

Hipótese	Variáveis	Teste	Estatística
Associação de 2 variáveis	Numéricas	Correlação	r
Associação de 2 variáveis	Categóricas	Qui-Quadrado	χ^2
Diferença de médias	1 numérica e 1 categórica	Teste t	t
Predição de valores	n variáveis de qualquer tipo e 1 alvo numérico	Regressão linear	b
Classificação	n variáveis de qualquer tipo e 1 alvo categórico	Regressão logística e classificadores	Odds-ratio etc

Erros

- Existem dois erros que queremos evitar no TSHN:
 - **Erro Tipo I (falso positivo):** dizer que um efeito existe quando ele não existe.
 - **Erro Tipo II (falso negativo):** dizer que um efeito não existe quando ele existe.
- Existem taxas máximas probabilísticas aceitáveis para esses erros:
 - Taxa de erro Tipo I: $\alpha = 0,05$ ou $\alpha = 0,01$
 - Taxa de erro Tipo II: $\beta = 0,20$ ou $\beta = 0,05$



Processo do TSHN

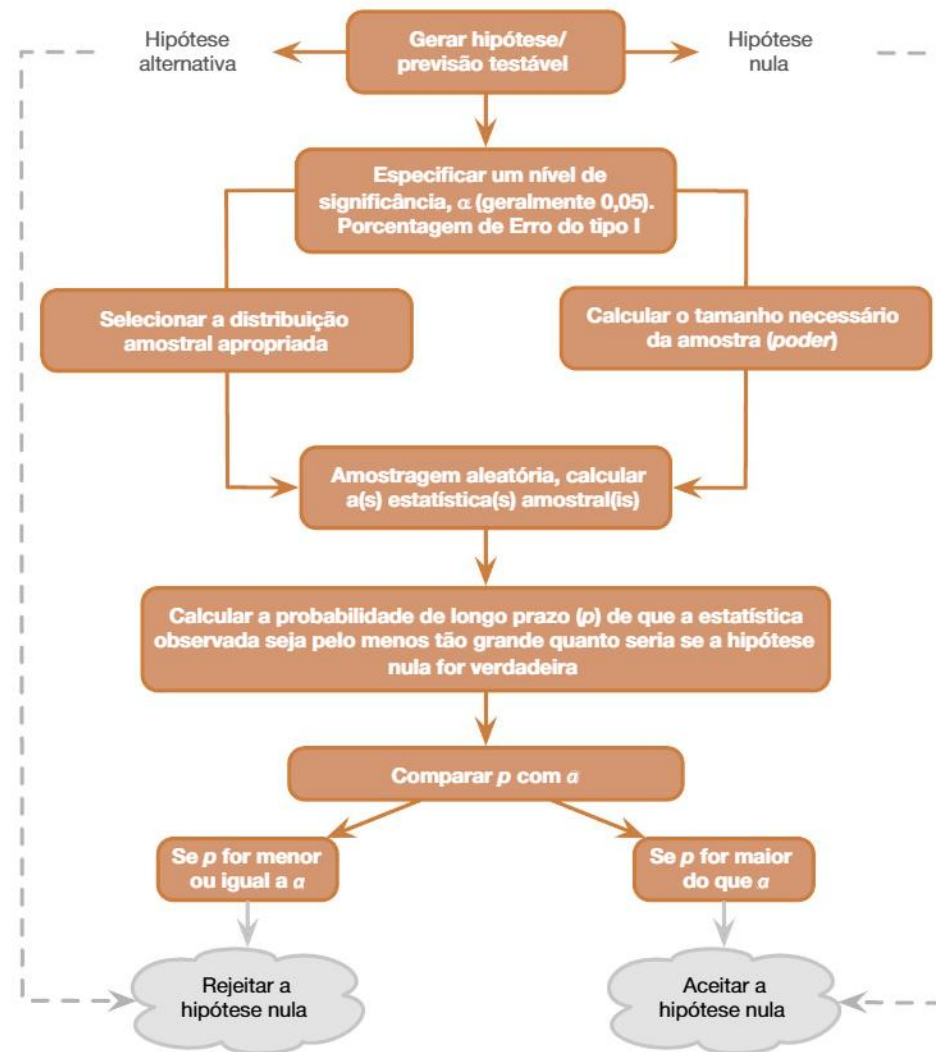


Figura 2.13 Fluxograma da testagem de hipóteses (TSHN).

Exemplo TSHN

- Vamos fazer um teste Z para verificar se estudantes brasileiros aceitos em universidades americanas tem escores significativamente mais elevados no TOEFL em comparação aos não aceitos.
-

Hipóteses

- h_0 : a diferença nos escores do TOEFL entre estudantes brasileiros que foram aceitos ou não em uma universidade americana não é estatisticamente significativa.
 - h_1 : a diferença nos escores do TOEFL entre estudantes brasileiros que foram aceitos ou não em uma universidade americana é estatisticamente significativa.
 - Vamos assumir que h_0 é verdadeira!
-

Alfa e distribuição amostral

- Vamos usar uma taxa de erro de tipo II de 0,05, o que significa que admitiremos uma chance de erro de 5 em 100, *i.e.*, em 5 de 100 estudos exatamente iguais a este estaremos errados!
 - Vamos supor que nossa estatística, uma diferença entre médias, vem de uma distribuição normal padrão, composta de escores Z.
-

Estatística

Grupo	Média	Variância	Diferença entre médias	Estatística Z
Aceito (n = 300)	102,47	25,87	102,47 - 110,32 = -7,85	-18,96
Não- aceito (n = 198)	110,32	17,07		

$$Z = \frac{\overline{x_1} - \overline{x_2}}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

Valor de p

- A estatística Z é um escore Z da DNP, logo, podemos usar a distribuição normal para saber a probabilidade dela. Usamos a distribuição normal padrão, com média = 0 e DP = 1 porque ela representa a hipótese nula, que supomos ser verdadeira.
- A probabilidade de um valor menor que ou igual a -18,96 é inferior a 0,001, assim, dizemos:

$$P(z \leq -18,96) < 0,001$$

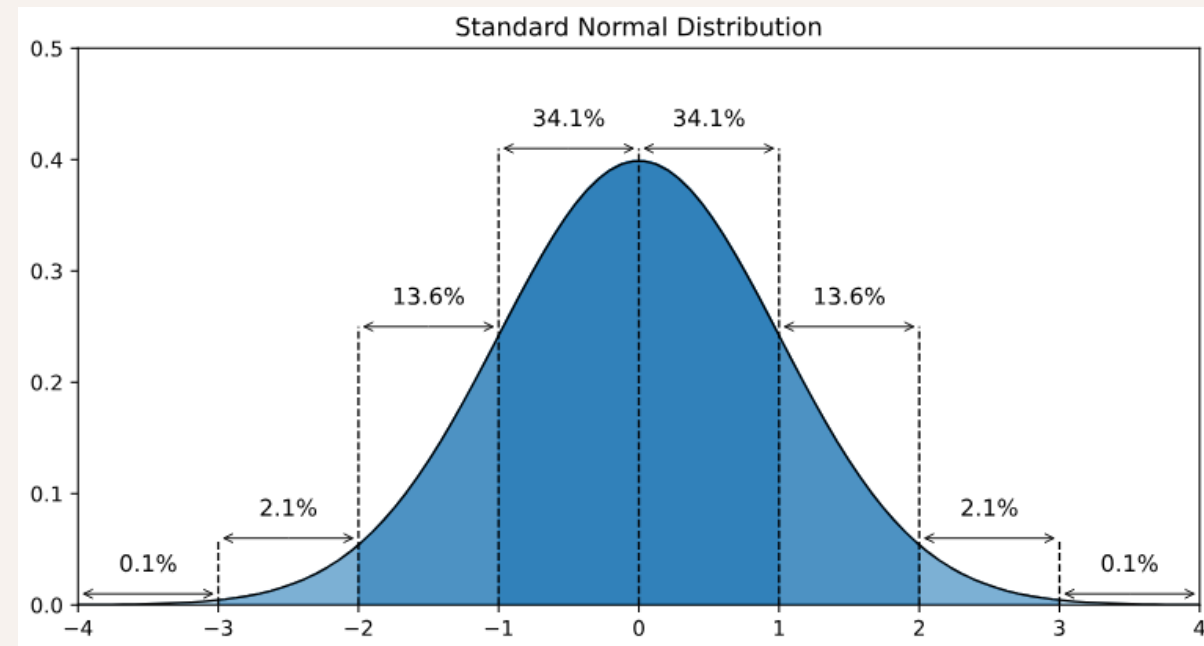
ou,

$$P(z \geq 18,96) < 0,001$$

Valor de p

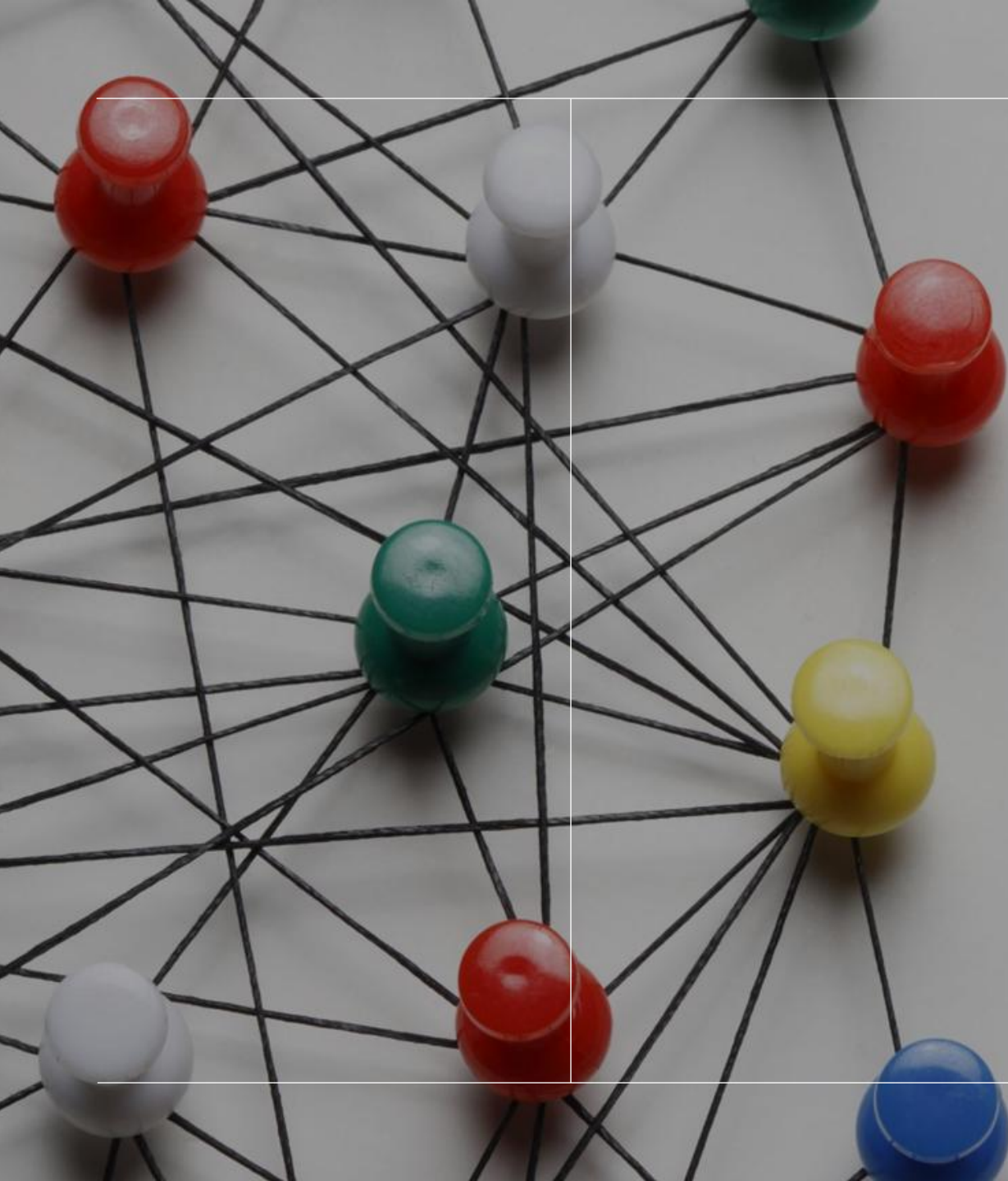
$$p < 0,001$$

- Significa que a probabilidade de encontrar uma estatística tão grande quanto essa ou maior é muito pequena.
- Ou seja, nosso resultado é muito improvável se supormos que a hipótese nula é verdadeira.



Interpretação

- Como nossa estatística é muito improvável sob a suposição de que a hipótese nula é verdadeira, rejeitamos a hipótese nula e assumimos a hipótese alternativa.
 - O TSHN funciona como um juri onde o suspeito é considerado inocente até que as evidências e a ideia de que ele seja inocente se torne ridícula.
-



bootstrap

Definição

- *Bootstrap* (alça de bota) é um termo usado para descrever uma expressão idiomática no inglês “*By one's own bootstraps*”, que sugere a tarefa impossível de alguém se levantar pelas alças das próprias botas. O termo se refere a um processo independente, autossustentável.
- Um procedimento de *bootstrap* ou reamostragem em estatística consiste em uma técnica de simulação que produz uma distribuição amostral.



Execução

1. Retirar uma amostra de tamanho n de um banco de dados com n casos, com reposição
 2. Calcular uma estatística de interesse
 3. Repetir os passos 1 e 2 k vezes ($k = 500$ ou $k = 1000$)
 4. Unir as k estatísticas calculadas em uma distribuição
 5. Usar a distribuição das estatísticas como uma distribuição amostral, calculando intervalos de confiança a partir dos percentis
-

Uso

- Quase todos os testes estatísticos podem ser usados com procedimentos de reamostragem.
 - Se os intervalos de confiança não cruzarem o 0, concluímos que as estatísticas são significativas.
 - Ver exemplo no JASP!
-

Para se aprofundar...

