

Introdução à Estatística

CCA5968-Métodos e técnicas aplicadas à pesquisa em Comunicação

Vinícius Alves Sarralheiro

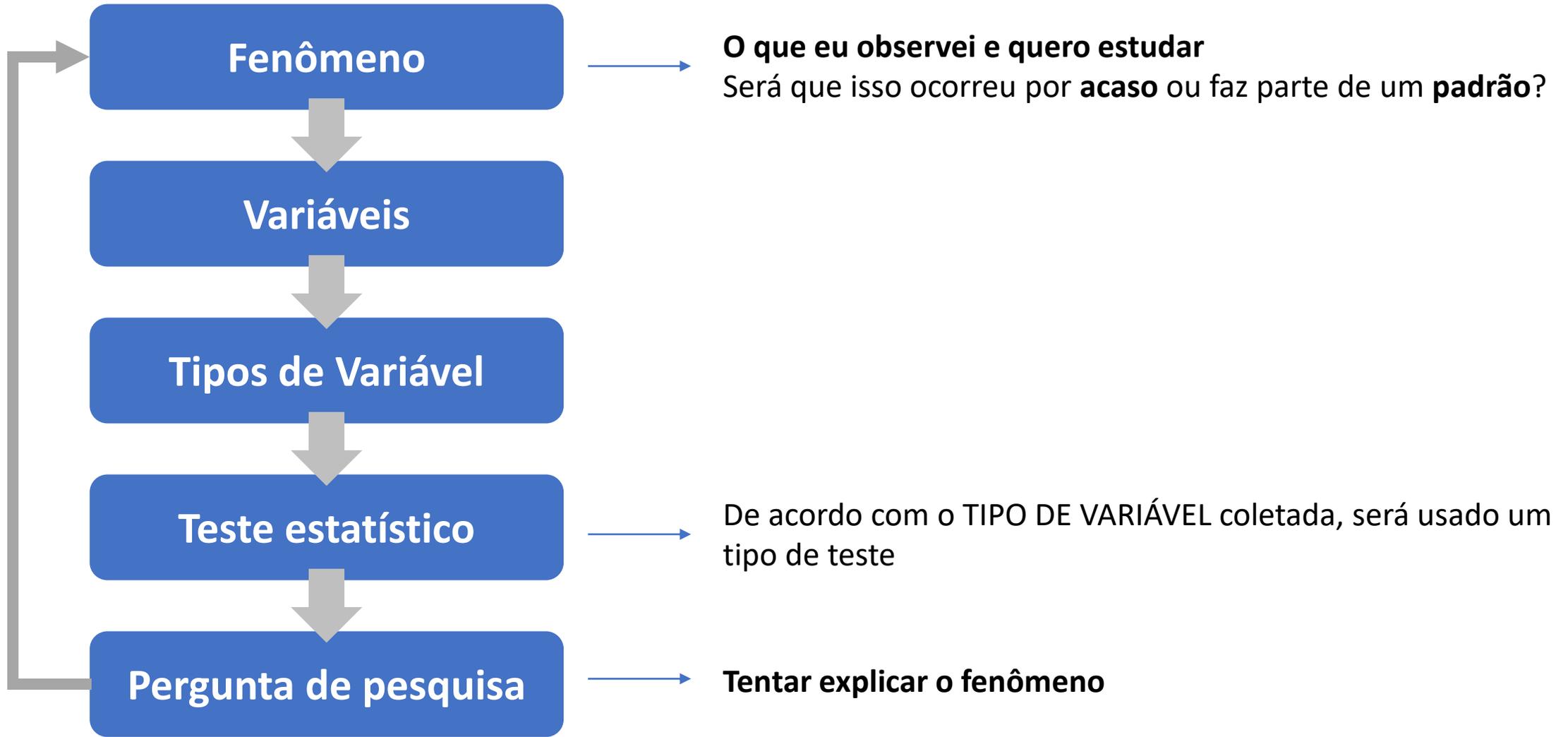


GENERALIZAÇÃO

Fenômeno → mensuração → salto interpretativo → generalização

ABRIR MÃO DAS CERTEZAS E DA VERDADE

Análise de Dados





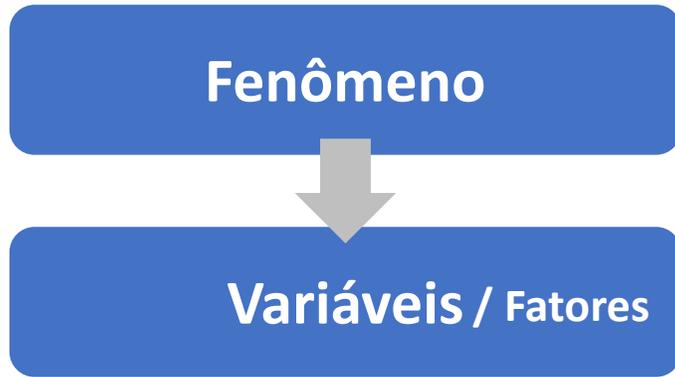
TEORIA DA MEDIDA

Em geral, criar réguas para a realidade



1) Validade: Será que o que eu estou medindo, de fato mede o que eu gostaria?

2) Precisão: Será que a “régua” que eu utilizo consegue discriminar coisas que tenho interesse?

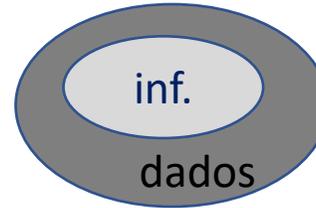


VARIÁVEL: Todo **dado** diretamente observável que pode ser medido

DADO → medida

≠ Informação → medida + sentido (ex: unidade de medida)

30 anos
 dado
 informação



Toda informação é um dado, mas nem todo dado é uma informação
SENTIDO / CONTEXTO

20 Ω

DIRETAMENTE OBSERVÁVEL → É percebida de modo direto, sem meios
 Todos podem reproduzir (é auditável)
 Ex: Peso, altura, nº do sapato, renda, cor preferida, temperatura...

Mas e quando não é diretamente observável?

Ex: satisfação, beleza, inteligência, racismo...

FATOR: Agrupamento de variáveis com o objetivo de descrever um atributo latente



7 = 7

Satisfação é um atributo interno



Dica: na dúvida, colete tudo o que puder de forma contínua

Ex: IMC, aprendizado (1 a 5)

a) **CATEGÓRICAS:** variável ou fator formada por categorias, onde cada uma expressa um atributo do fenômeno

QUALITATIVAS

Ex: Time de futebol, gênero, cor do cabelo, cor preferida

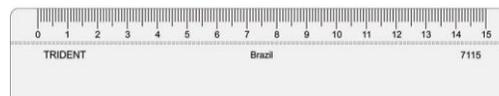


- 1) Nominal (nomes sem ordem)
- 2) Ordinal (categorias ordenadas)
- 3) Binárias (2 categorias)

b) **CONTÍNUAS:** variável ou fator que possui uma forma de apresentação linear, com infinitos níveis de observação

QUANTITATIVAS

Ex: idade, peso, altura



- 1) Contínuas
- 2) Discreta (ex. escala Likert)

Fenômeno



Variáveis



Tipos de Variável



Teste estatístico



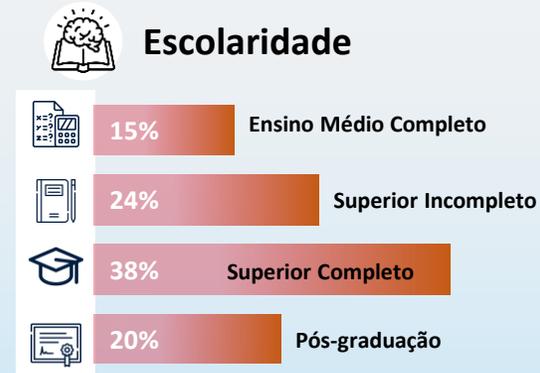
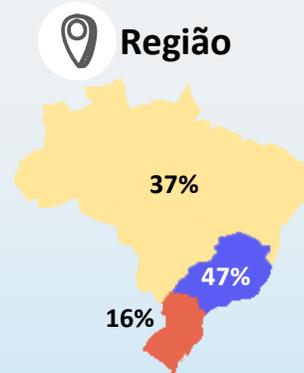
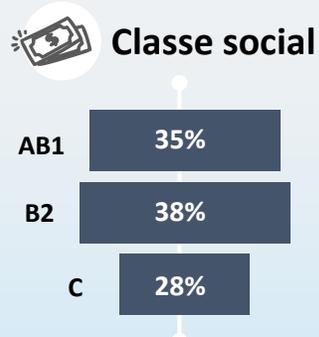
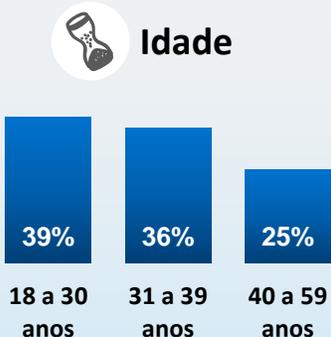
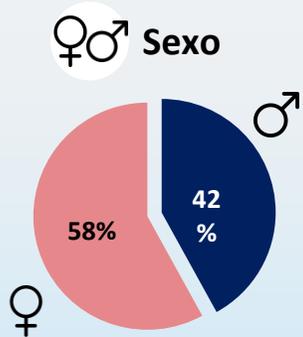
De acordo com o TIPO DE VARIÁVEL coletada, será usado um tipo de teste

MEDIDAS DESCRITIVAS (FORMAS DE MENSURAÇÃO):
Resumos iniciais para descrever os dados

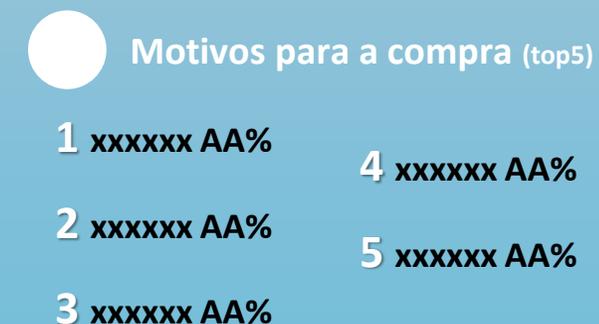
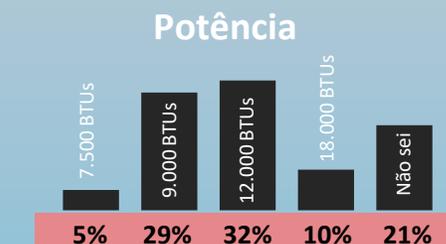
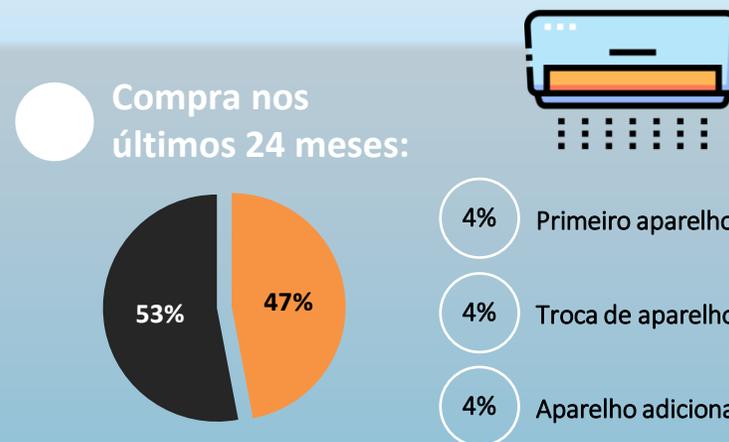
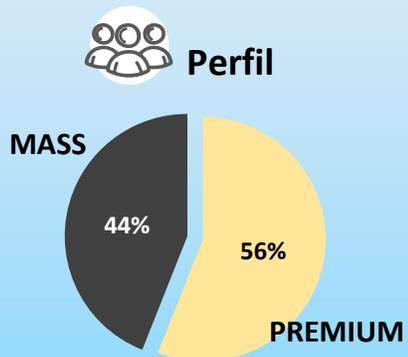


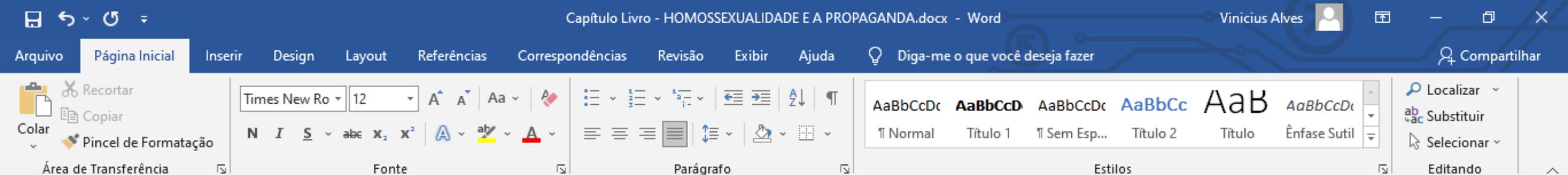
PERFIL

Medidas descritivas (exemplo)



55% Não tem filhos





estímulo, de modo que cada participante visualizasse uma sequência diferente de publicações.

Os dados foram coletados no Laboratório 4C¹ da ECA-USP, utilizando aparelho rastreador de olhos Tobii X2-60, posicionado na base de um monitor LCD de 19 polegadas que, por sua vez, estava a uma distância de 65cm dos participantes. As condições de luminosidade do ambiente foram controladas e mantidas estáveis. Para a análise dos dados foi utilizado o software Tobii Studio versão 3.4.2.

Procedimentos metodológicos

O estudo foi realizado com 30 sujeitos, de 18 a 32 anos, declarados como pertencentes ao sexo e ao gênero masculinos (portanto, homens cisgênero), sendo metade (15) heterossexuais e metade (15) homossexuais, todos por auto-declaração. Os participantes deveriam se encaixar em um desses grupos, que se tornarão a base para a análise, e a seleção para compor cada um deles foi por conveniência e aleatória, sendo sua maioria formada por alunos da Universidade de São Paulo. Na entrevista de recrutamento não se notou viés algum em termos de atitudes relevantes para o projeto. Os participantes abordados eram apresentados aos objetivos e, após o consentimento individual, eram informados de que o *eye tracker* era o instrumento utilizado para a coleta dos dados. Visando ao anonimato, os nomes dos participantes foram codificados utilizando um sistema alfanumérico (ex. P01).

O aparelho foi calibrado individualmente e iniciada a apresentação dos estímulos. A apresentação dos anúncios ocorreu de forma aleatória, sendo que cada imagem ficou em exposição por aproximadamente 10 segundos, tempo ideal definido pelo pré-teste e possivelmente similar ao despendido em uma leitura de revista. Ao final da exibição de cada anúncio, uma pergunta era feita a cada participante no próprio aparelho de coleta de dados a fim de abordar a atitude em relação ao anúncio: “Quanto você gosta desse anúncio?”; a questão era acompanhada por uma escala de notas que iam de 0 a 10 e a escolha dos critérios para atribuição dessa nota era livre para cada um dos sujeitos.

Fenômeno



Variáveis



Tipos de Variável



Teste estatístico

De acordo com o TIPO DE VARIÁVEL coletada, será usado um tipo de teste

MEDIDAS DESCRITIVAS (FORMAS DE MENSURAÇÃO):

Resumos iniciais para descrever os dados

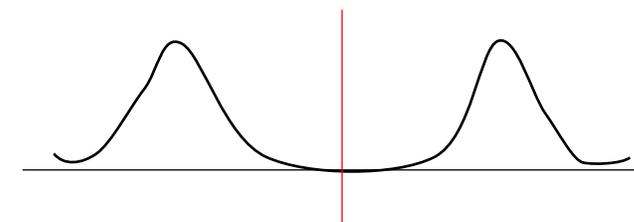
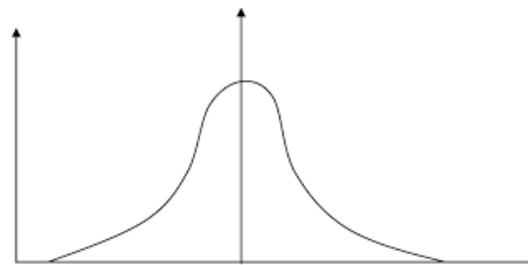
1) CONTÍNUAS: Média, Variância, Desvio Padrão, Mediana, Moda

Peso (kg)
3
2
1

$$\text{Média} = \frac{\text{Soma das obs}}{\text{N obs}} = \frac{3+2+1}{3} = 2 \text{ kg}$$

Média = estimativa do valor esperado → palpite

Medida de centralidade – ideia de palpite onde os dados estão em média
Reduz a incerteza do “chute”, porém não tenho informação da distribuição dos dados



*Preciso de uma medida de dispersão

MEDIDAS DESCRITIVAS (FORMAS DE MENSURAÇÃO):

Resumos iniciais para descrever os dados

1) CONTÍNUAS: Média, Variância, Desvio Padrão, Mediana, Moda

Média = 2kg

Peso (kg)	Desvio	Média (desvio) ²
3	3-2 = 1	1
2	2-2 = 0	0
1	1-2 = -1	1

O quanto os dados se distanciam da média, em média

Tira o efeito do sinal

Problema de comparação

$$\text{Variância} = \frac{\sum (\text{desvios})^2}{N \text{ obs} - 1} = \frac{1+0+1}{2} = 1 \text{ kg}^2$$

Valor ind. - média



Medida de **distância** individual em relação à média

Desvio Padrão → mesma definição da Variância (O quanto os dados se distanciam da média, em média)

$$DP = \sqrt{\text{variância}}$$

Mesma unidade da média → DP = 1 kg

Média e DP = INTERVALO PARA O MEU PALPITE



MEDIDAS DESCRITIVAS (FORMAS DE MENSURAÇÃO):

Resumos iniciais para descrever os dados

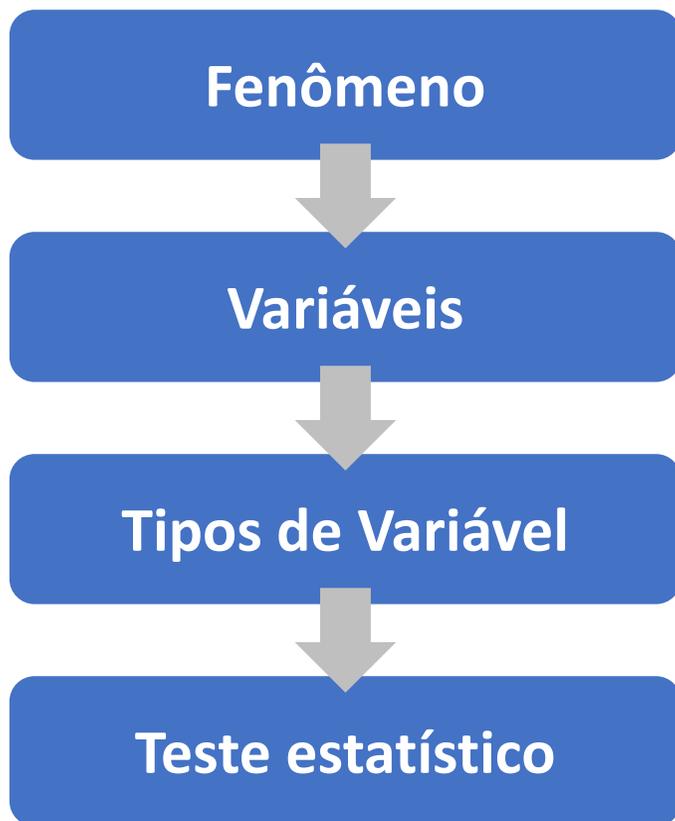
1) CATEGÓRICAS:

a) Número de observações (N) → Medida absoluta
Ex: Do total de 100 na amostra, são 60 homens e 40 mulheres

b) Porcentagem (%) → Número relativo
60% homens e 40% mulheres

c) Moda: Categoria mais frequente

d) Mediana: Categoria na centralidade
(Somente para categóricas ordinais, pois depende de uma ordem)



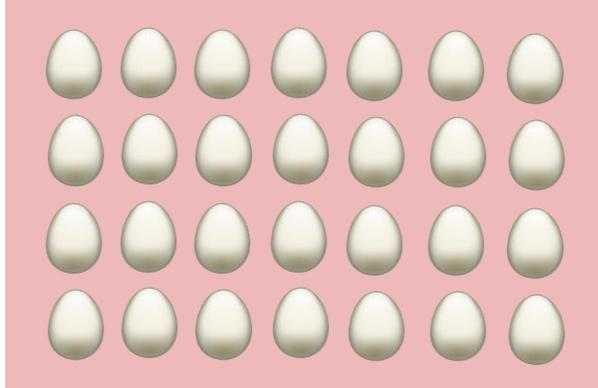
TESTES ESTATÍSTICOS

1) Efeito de variáveis categóricas sobre variáveis contínuas:
Teste t, Anova, Manova, GLM

2) Associação entre variáveis categóricas:
Qui-quadrado, Análise de Correspondência

3) Fatores associados com uma variável categórica OU Fatores preditivos para uma variável contínua:
Correlação, **Análise Fatorial**, Regressão Linear, Regressão Logística

AMOSTRA



Pergunta: Será que todos os ovos são de galinha?

1) Pegar um ovo de cada vez e verificar se ele é de galinha

2) Abrir a caixa, dou uma olhada e se um deles **NÃO** for de galinha, respondi minha pergunta

Mas se todos, nessa olhada, forem de galinha, não significa que todos sejam de galinha

Procurar diferenças

Ovos \neq ovos de galinha



As pessoas votam no candidato X?



As pessoas compram o produto Y?

O pesquisador sempre busca encontrar DIFERENÇAS

Métodos de amostragem

1) Amostra Aleatória Simples (AAS)

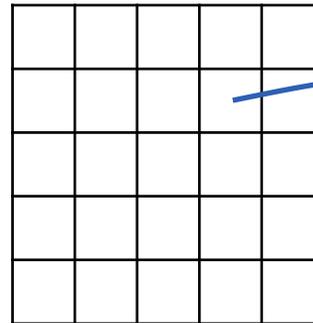
Processo em que todas as unidades de análise têm a mesma chance de serem selecionadas

SORTEIO → Vieses de azar e da representatividade da população (ex: sair todos homens)

2) Estratificado (estratos/quotas)

Pela característica da unidade de análise

Sexo (2)	}	48
Idade (6)		
Renda (4)		



AAS

Sorteio na população até cumprir as quotas

3) Conglomerado (cluster)

A premissa é, em geral, a questão geográfica

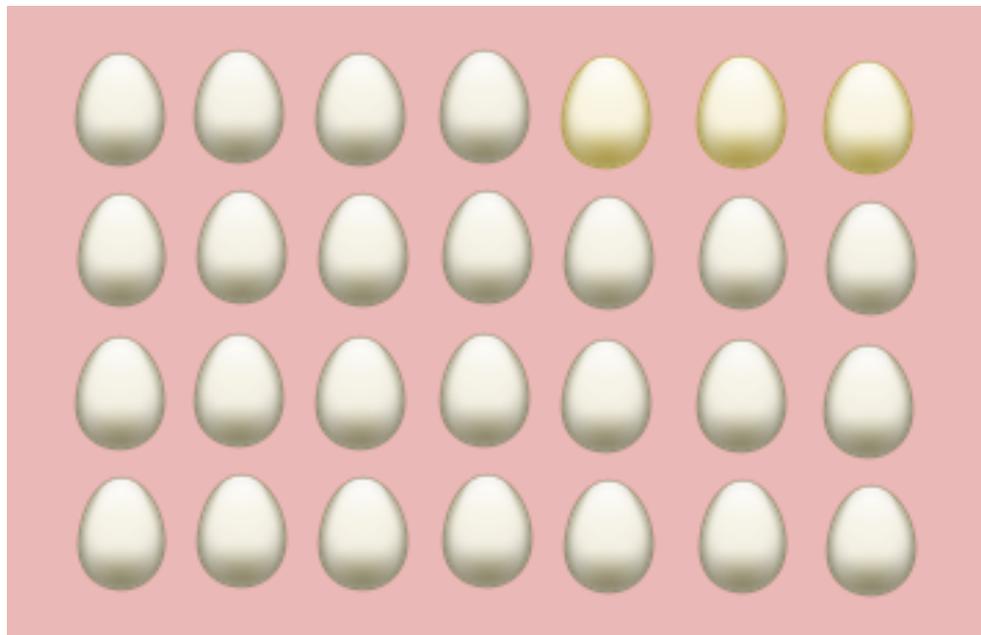
O indivíduo deve ter o perfil que represente sua região



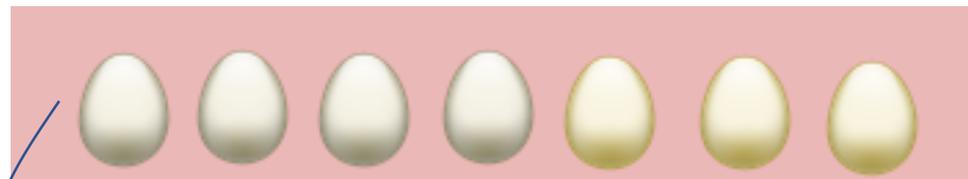
4) Amostra por Conveniência

É o que tem para hoje! → Se preocupar com os vieses que atrapalham as inferências

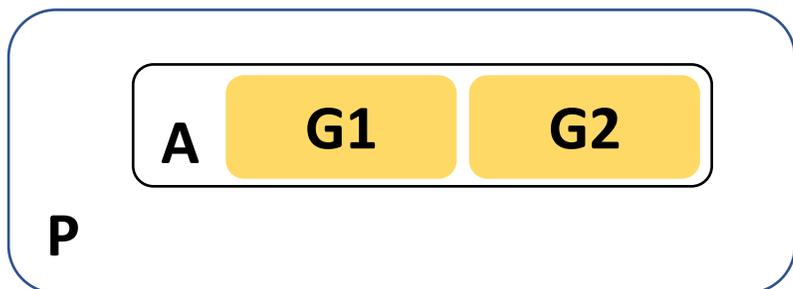
POPULAÇÃO



AMOSTRA



GRUPOS



O pesquisador sempre busca encontrar DIFERENÇAS

Partindo dessa ideia, cria-se a lógica científica

TESTE DE HIPÓTESES

O teste de hipóteses é uma suposição ou afirmação relativa a uma ou mais amostras, baseada nos parâmetros populacionais, que pode ser verdadeira ou falsa.

A ideia básica é que a partir de uma amostra da população será estabelecida uma regra de decisão, segundo a qual a hipótese proposta será não rejeitada ou rejeitada. **Esta regra de decisão é chamada de teste.**

Então, buscamos encontrar diferenças existentes entre os grupos na amostra estudada **Esse é o princípio básico do teste estatístico!**

O “p” (*p-value*)

Teste estatístico

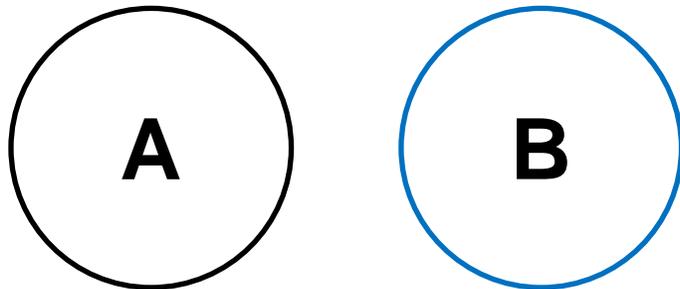
No fim, o que quero saber:

O quão certo ou errado estou em afirmar a hipótese (H_a)? $\rightarrow p$

Definição “inicial” do p:

Chance de erro em afirmar que coisas são diferentes (H_a) \rightarrow tamanho da evidência

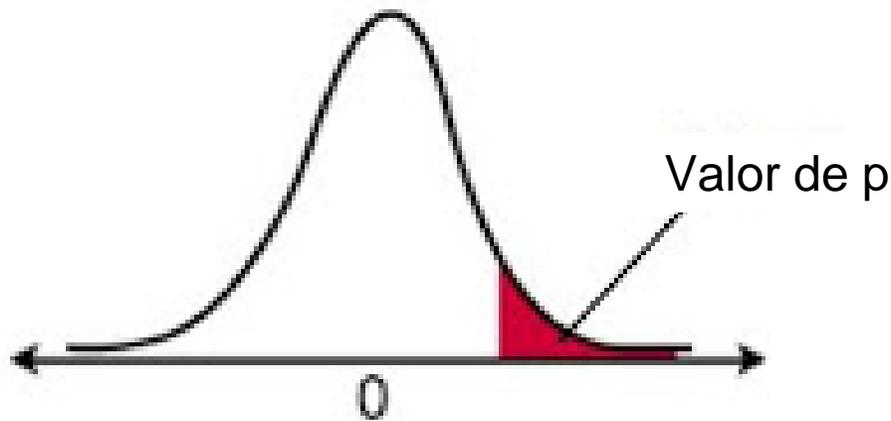
Chance de erro que eu aceito assumir em minha pesquisa



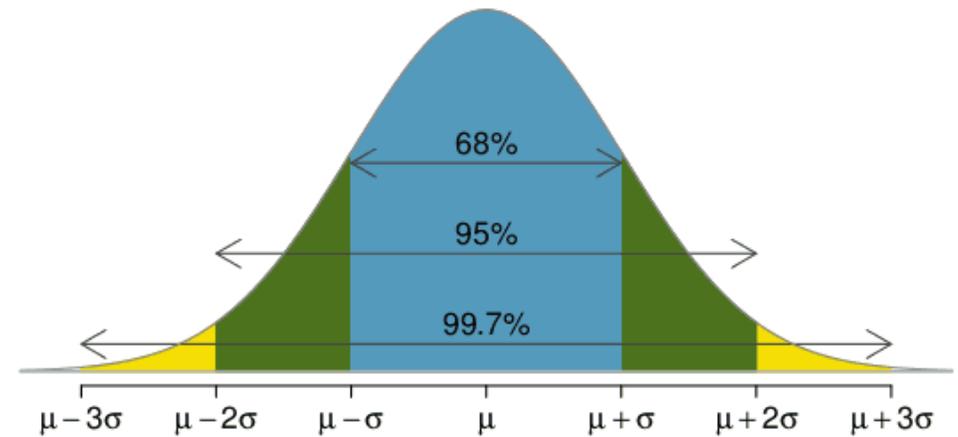
Posso afirmar que o grupo A é \neq do grupo B?

Quanto maior a diferença entre os grupos, menos errado eu estou
Portanto, quanto maior a diferença, menor é o p

O teste estatístico aplicado, vai ter dar um valor de “p” correspondente



Regra do 68, 95, 99

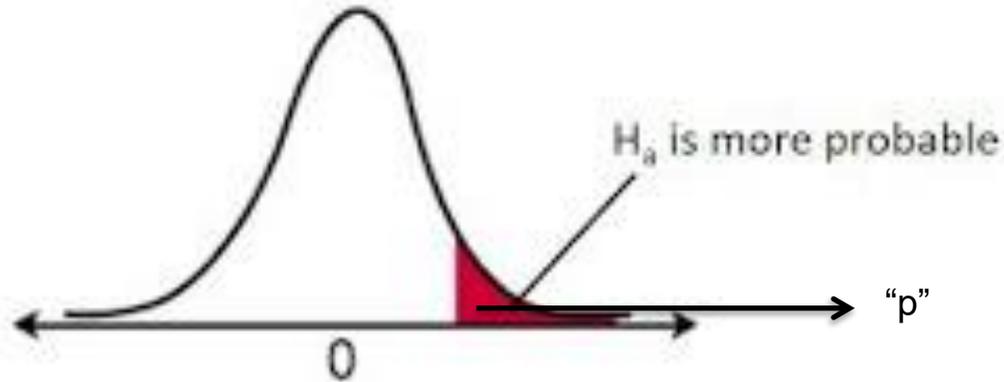


Valor padrão de “p” adotado na nossa área: **$p \leq 0,05$**

Definição formal do “p”:

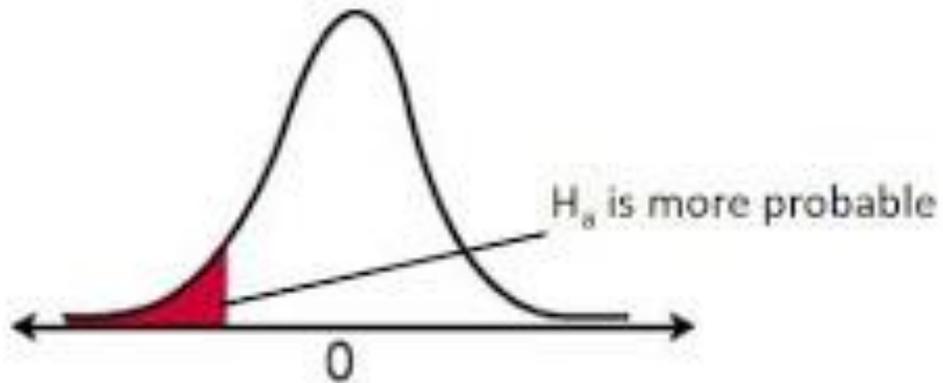
Probabilidade de encontrar uma estatística de interesse mais extrema do que a encontrada

Types of Hypothesis Tests



Right-tail test

$$H_a: \mu > \text{value}$$



Left-tail test

$$H_a: \mu < \text{value}$$

Se "p" > 0.05 → H₀ (a favor da =)

→ Muita chance de erro ao afirmar a diferença

Se "p" < ou = 0.05 → H_a (a favor da ≠)

→ Pouca chance de erro ao afirmar a diferença

RESUMINDO... CONSTRUÇÃO DE UM TESTE DE HIPÓTESES

- 1) Formulação da hipótese (H_a) a ser testada;
- 2) Use a teoria estatística e as informações disponíveis para decidir qual teste (estimador) será usado para testar a hipótese H_a ;
- 3) Use as observações da amostra para calcular o valor da estatística do teste (Medidas Descritivas);
- 4) Verifique o “p” encontrado para poder afirmar se existe ou não diferença no fenômeno testado;
- 5) Se “p” $>$ 0.05 (5%), muita chance de erro em afirmar que existe diferença
Se “p” $<$ ou $=$ 0.05 (5%), aceitamos a H_a , pois há pouca chance de erro em afirmar a diferença

Teste t

Comparar DOIS grupos independentes em função de suas médias

PERGUNTA: Existe efeito de uma variável categórica com dois grupos sobre uma variável contínua?

GRUPO	VARIÁVEL
A	-----
B	-----

cat cont

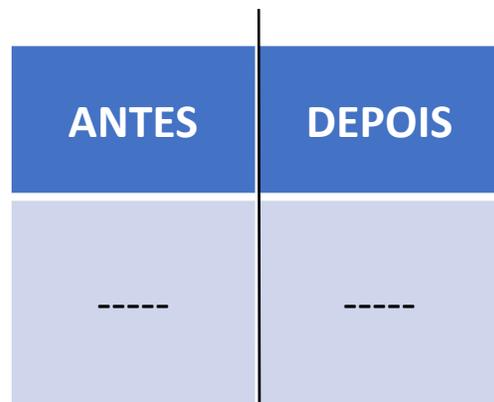
H0: Méd A = Méd B
HA: Méd A ≠ Méd B

$$t = \frac{\text{Média A} - \text{Média B}}{Dp (A - B)}$$

Teste t

Varição ao longo do tempo de duas médias (pareado / medidas repetidas)

PERGUNTA: Existe efeito do [tempo, comercial, treinamento...] sobre uma variável contínua?



[tempo, comercial, treinamento...]

H0: Méd $\Delta = 0$
HA: Méd $\Delta \neq 0$

$$t = \frac{\text{Média } \Delta}{\frac{\text{Dp } \Delta}{\sqrt{n} \text{ amostral}}}$$

Média ANTES – Média DEPOIS = Δ

Grupo	Intenção de compra (escala de 0 a 10)			Avaliação do comercial (notas 0 a 10)		
	N	Média	Desvio Padrão	N	Média	Desvio Padrão
Viu Comercial 1	500	9,02	1,46	500	7,13	2,67
Viu Comercial 2	500	8,72	1,47	500	7,20	2,59
Total	1000	8,93	1,47	1000	7,15	2,65

Teste t

		Statistic	df	p
Intenção de compra	Student's t	3.125	999	0.002
Avaliação (nota)	Student's t	-0.379	999	0.705

Qui-quadrado

PERGUNTA: Existe associação entre 2 ou mais variáveis categóricas?

Como não sei o que esperar, começo pela frequência observada

TABELA DE CONTINGÊNCIA

GRUPO	ANÚNCIO
A	1
B	2

	1	2	
A	3	2	→ 5
B	6	10	→ 16
	↓ 9	↓ 12	21

	1	2
A	$\frac{5 \cdot 9}{21} = 2$	$\frac{12 \cdot 5}{21} = 3$
B	$\frac{9 \cdot 16}{21} = 7$	$\frac{12 \cdot 16}{21} = 9$

$$\chi^2 = 0,5 + 0,3 + 0,1 + 0,1$$

$$\chi^2 = 1,0$$

$$\chi^2 = \frac{(\text{RESÍDUO})^2}{\text{ESPERADO}}$$

Frequência esperada:
 $\sum \text{coluna} \cdot \sum \text{linha} / \text{total}$

$$\chi^2 = (-1)^2/2 + (1)^2/3 + (1)^2/7 + (-1)^2/9$$

Rede social mais utilizada pelo participante	Conhece o meme?			Total
		Sim	Não	
Facebook	% entre linhas	50,5%	49,5%	100%
Pinterest	% entre linhas	18,3%	81,7%	100%
Twitter	% entre linhas	71,4%	28,6%	100%
TikTok	% entre linhas	42,1%	57,9%	100%
Instagram	% entre linhas	63,4%	36,6%	100%
Total	% entre linhas	49,1%	50,9%	100%

Qui-quadrado (χ^2)	Value	df	p
	19.3	4	< .001
N	1167		

Pergunta de pesquisa: Propagandas de produtos de limpeza protagonizadas por mulheres são mais bem recebidas pelo público?

QUI-QUADRADO

PERGUNTANDO QUAL A PROPAGANDA PREFERIDA:

PERGUNTA: Existe associação entre 2 variáveis categóricas?

		Tipo de propaganda preferida	
		c/ Mulheres	c/ Homens
Sexo	Masc	X	Z
	Fem	Y	W

$$\chi^2 = \frac{(\text{RESÍDUO})^2}{\text{ESPERADO}}$$

TESTE t

PERGUNTANDO A NOTA DE CADA PROPAGANDA

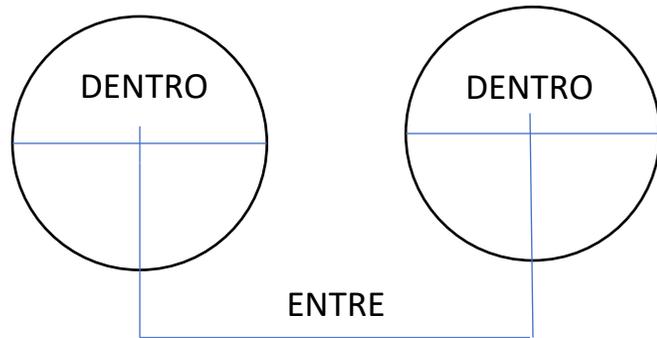
PERGUNTA: Existe efeito de uma variável categórica com dois grupos sobre uma variável contínua?

Sexo	Nota Propag. c/ Mulheres	Sexo	Nota Propag. c/ Homens
Masc	0 a 10	Masc	0 a 10
Fem	0 a 10	Fem	0 a 10

$$t = \frac{\text{Média Masc} - \text{Média Fem}}{Dp (\text{Masc} - \text{Fem})}$$

ANOVA

PERGUNTA: Existe efeito de uma variável categórica com dois ou mais grupos sobre uma contínua?



variabilidade ENTRE

$$\frac{N \text{ do grupo} * (\text{Média grupo} - \text{Média geral})^2}{n^{\circ} \text{ de grupos} - 1}$$

variabilidade DENTRO

$$\frac{(\text{valor} - \text{Média grupo})^2}{N \text{ amostral} - n^{\circ} \text{ de grupos}}$$

$$F = \frac{\text{variabilidade ENTRE}}{\text{variabilidade DENTRO}} = \frac{N \text{ do grupo} * (\text{Média grupo} - \text{Média geral})^2}{n^{\circ} \text{ de grupos} - 1} \div \frac{(\text{valor} - \text{Média grupo})^2}{N \text{ amostral} - n^{\circ} \text{ de grupos}}$$

Dois (ou mais) grupos são diferentes quando a variabilidade ENTRE os grupos é grande e DENTRO dos grupos é pequena

	Intenção de compra (escala de 0 a 10)			Avaliação do comercial (notas 0 a 10)		
Grupo	N	Média	Desvio Padrão	N	Média	Desvio Padrão
Viu Comercial 1	500	9,02	1,46	500	7,13	2,67
Viu Comercial 2	500	8,72	1,47	500	7,20	2,59
Total	1000	8,93	1,47	1000	7,15	2,65

ANOVA

	Sum of Squares	df	Mean Square	F	p
Intenção de compra	20.8	1	20.84	9.77	0.002
Avaliação (nota)	1.01	1	1.01	0.144	0.705

Análise Fatorial

O objetivo é criar fatores (medidas latentes)

Usos principais:

- 1) Redução do número de variáveis (bloco) → Fator conjunto de variáveis para formar um bloco
- 2) Validação de questionário ou escala → encontrar o mesmo agrupamento original
- 3) Criação de índices → fator único

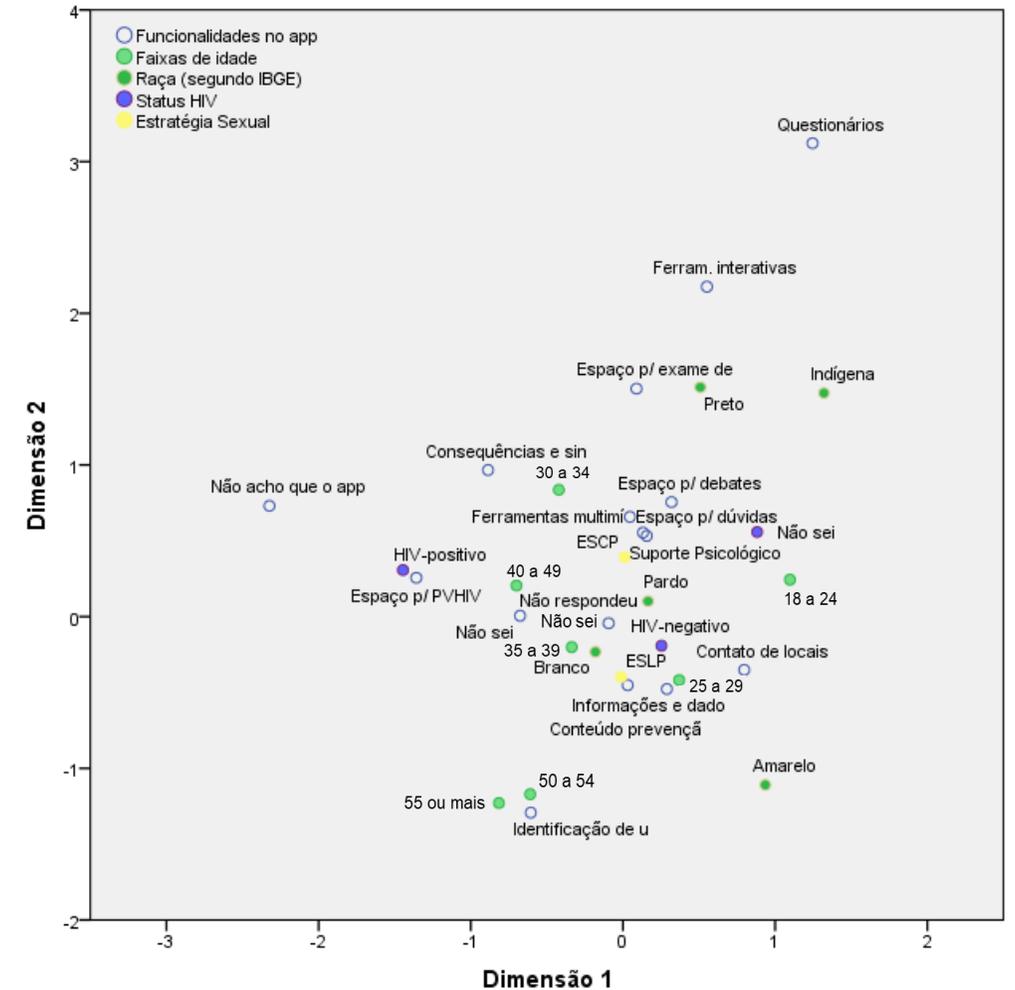
Variáveis	FATORES
V1	F1
V2	
V3	
V4	F2
V5	
V6	F3
V7	

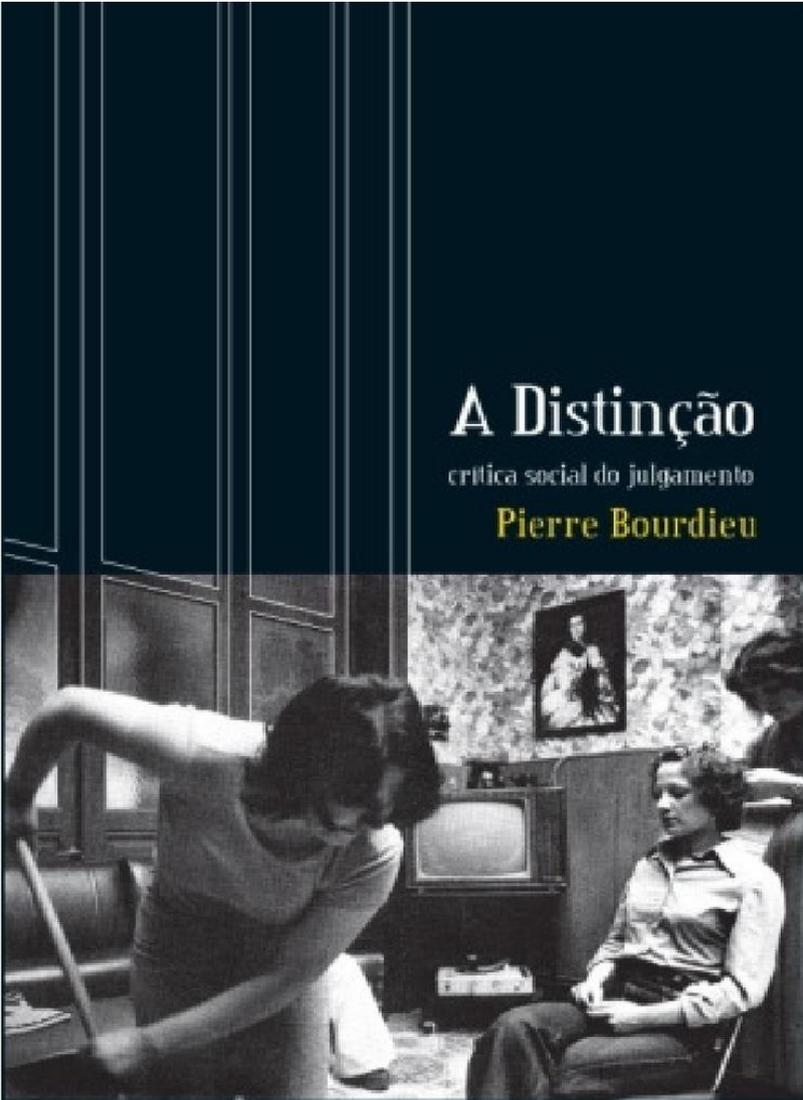
Análise de Correspondência

Gerar perfis a partir da associação entre variáveis categóricas

OLHAR “QUALITATIVO” SOBRE OS DADOS

- Formar perfis a partir das variáveis
- Segue a mesma lógica de associação do qui-quadrado, mas de forma múltipla
- Considera as dimensões, que devem ser interpretadas também de forma qualitativa





O gosto dominante

A boa vontade cultural

O gosto de necessidade

Quali x Quanti?



Análise textual (CAMARGO; JUSTO, 2013)

Organização dos textos → Unidades de significação →
Categorização dos significados → Meta-textos de análise

- 1) Organização e preparo dos dados para a análise;
- 2) Transcrição e leitura de todos os dados, com releituras para avaliação do conteúdo transcrito e organização das unidades de significado;
- 3) Análise do material utilizando o software IRaMuTeQ, a partir de três análises principais – análise lexicográfica, análise de Classificação Hierárquica Descendente (CHD) e Análise Fatorial de Correspondência (AFC);
- 4) Utilização do processo de codificação para descrever o cenário, com avaliação das classes e nova escuta das entrevistas;
- 5) Descrição e apresentação dos temas e narrativas encontrados, sustentados pela literatura.

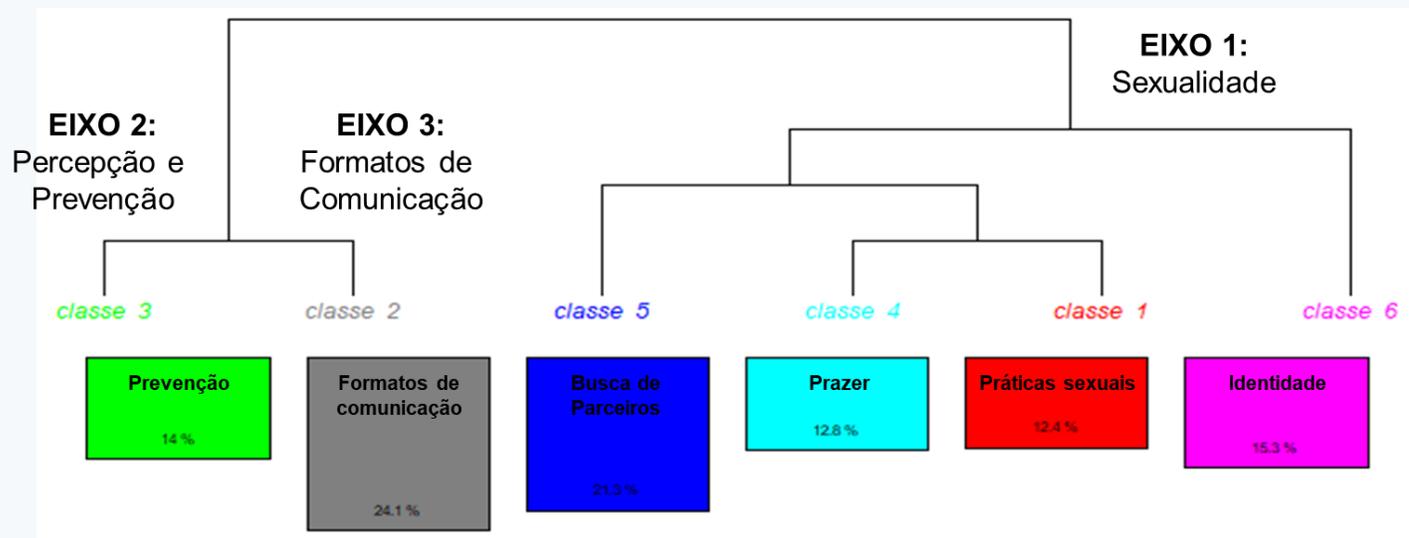
Parâmetros de qualidade (Mendes et al, 2019):

- a) aproveitamento mínimo de 70% dos STs na análise lexicográfica;
- b) valor de qui-quadrado superior a 3,84 (e, portanto, $p < 0,0001$) nas análises de CHD, (proximidade entre as palavras e a separação entre classes satisfatórias);
- c) soma dos fatores da AFC mais próxima a 100%.



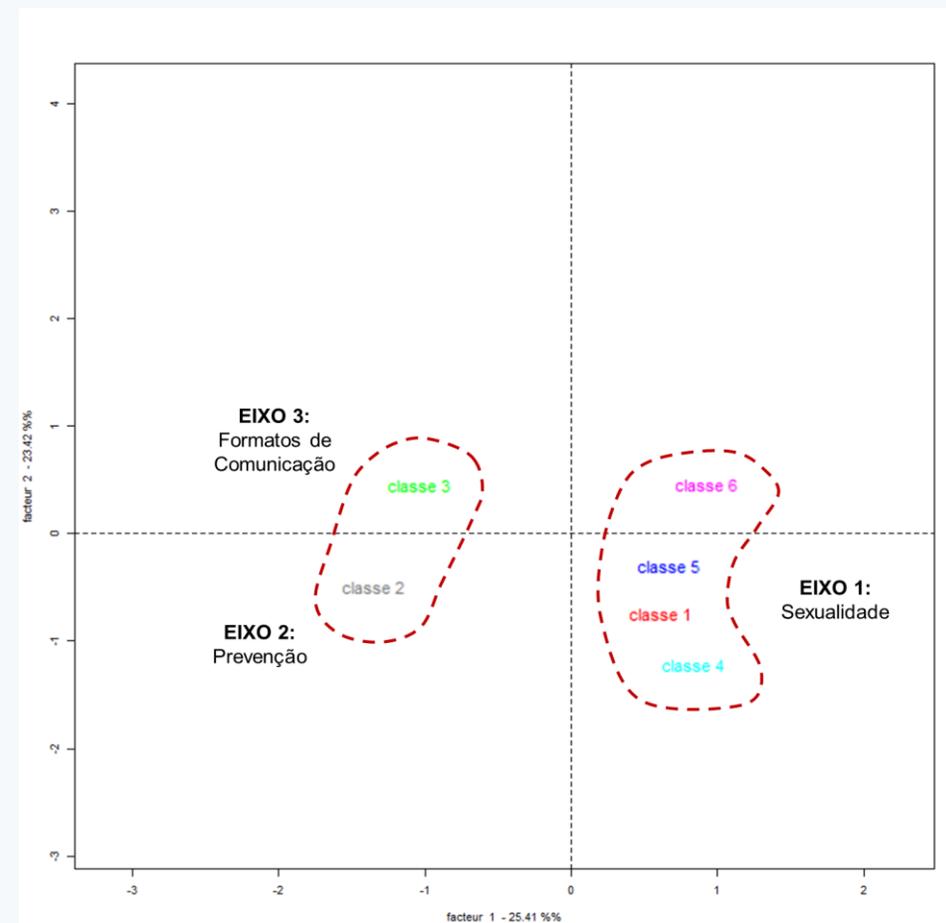
<http://www.iramuteq.org/>





aids	informação	aplicativo	sentir	camisinha	gay
HIV	campanha	Hornet	gozar	sexo	homossexual
exame	governo	Grindr	tesão	oral	homem
PrEP	carnaval	cara	prazer	penetração	família
doença	propaganda	conversa	vontade	ativo	mulher
pegar	medo	foto	momento	passivo	íntimo
sífilis	saúde	perfil	gostoso	pau	afetivo
proteção	comunicação	abordagem	confortável	rola	Deus
soropositivo	vídeo	mensagem	confiança	contato	corpo
DST	estigma	encontrar	motel	tentar	sexualidade
prevenção	pânico	Twitter	fetich	beijo	amar
contrair	população	Instagram	drogas	parar	parecer
PEP	atingir	mandar	cama	casa	caráter
remédio	lembrar	interesse	risco	público	assumir
vírus	redes sociais	relacionamento	gostar	sensibilidade	humano

Dendrograma das classes lexicais obtidas a partir da classificação hierárquica descendente (CHD) das palavras ativas advindas das entrevistas em profundidade



Análise Fatorial de Correspondência (AFC) das classes lexicais obtidas na Classificação Hierárquica Descendente (CDH) das entrevistas em profundidade

Teste t:

- 1) Comparar DOIS grupos independentes em função de suas médias
Uma variável categórica (2 grupos) e uma variável contínua
- 2) Variação ao longo do tempo de DUAS médias (pareado / medidas repetidas)
Antes e depois como variável

ANOVA

- 1) Comparar DOIS OU MAIS grupos independentes
Uma variável categórica e uma ou mais variáveis contínuas
- 2) Variação ao longo do tempo de DUAS OU MAIS médias (pareado / medidas repetidas)
Efeito 1, Efeito 2, Efeito 3...

QUI- QUADRADO

Associação entre DUAS OU MAIS variáveis
Somente para variáveis categóricas

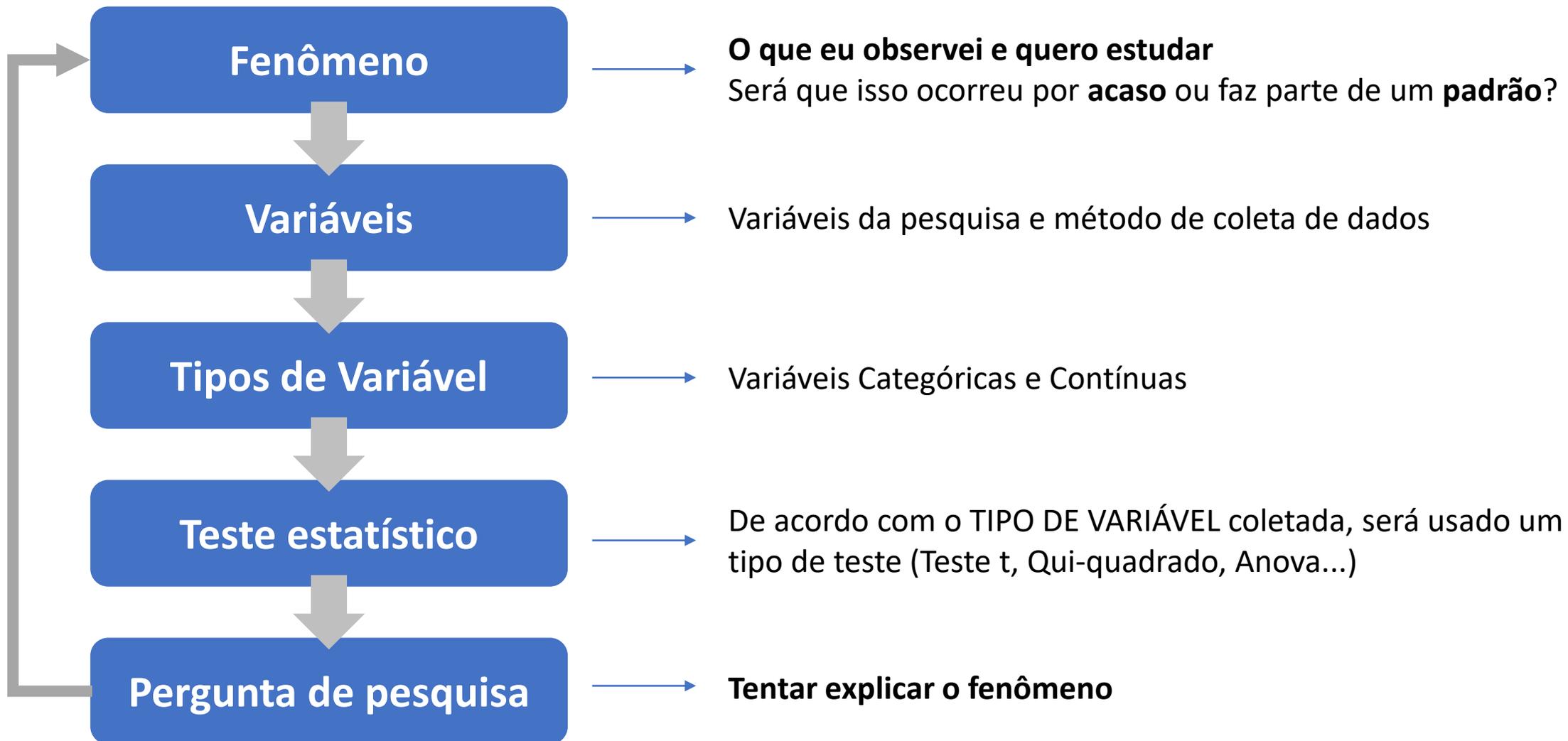
ANÁLISE FATORIAL

Criar medidas latentes (fatores) → escalas (contínuas)
Análise Fatorial Exploratória e Análise Fatorial Confirmatória

ANACOR

Gerar perfis a partir da associação entre variáveis categóricas

Análise de Dados



Análise de Dados

