

Estatística

Conceitos

Estatística Descritiva: apresenta métricas sobre os dados.

Estatística Inferencial: realiza previsões a partir dos dados existentes.

C/qual conjunto será trabalhado?	
Conj. Completo	Conj. Incompleto
População	Amostra
Censo	Amostragem
Parâmetro	Estimativa

Tipos de variáveis			
Quantitativas		Qualitativa	
Discretas	Continuas	Nominais	Ordinais
valores inteiros	valores inteiros e fracionários	s/ordem entre valores	c/ordem entre valores
Gtd. objetos	Altura, peso.	Gênero, cor, país	Grau de escolaridade

Métricas		
Elemento	População	Amostra
	Parâmetro populacional θ	Estimadores do par. populacional $\hat{\theta}$
Média	μ	\bar{x}
Variância	σ^2	s^2
Desvio Padrão	σ	s
Tamanho	N	n
Proporção	P	\hat{p}
Propriedades dos estimadores: - não viesado: na média, acerta o parâmetro populacional - eficiente: não viesado e o mais preciso possível (menor variância possível) - consistente: à medida que a amostra cresce, será convergido para valor do parâmetro - Máxima verossimilhança: estimador mantém a mesma distribuição de probabilidade da população		

Estatística Descritiva

1. Formas de apresentação dos dados

☐ **Dados brutos:** listagem de todos os dados
 $\{0,5, 10, 15, 15, 15, 20, 20, 30\}$

☐ **Dados ponderados:** tabela de frequência sem intervalo

* n de classes: amplitude/n

* n de intervalos: \sqrt{n}

Valor Observado (Xi)	Frequência Absoluta (fi)	Frequência Relativa (fri)	Frequência Acumulada (Fi)	Frequência Acumulada Relativa (Fri)
0 —10	2	2/9 \cong 22%	2	2/9 \cong 22%
10 —20	4	4/9 \cong 44%	6	6/9 \cong 67%
20 —30	3	3/9 \cong 33%	9	9/9 = 100%
Soma (Σ_i)	9	9/9 = 100%	-	-

☐ **Dados agrupados:** tabela de frequência com intervalo

Valor Observado (Xi)	Frequência Absoluta (fi)	Frequência Relativa (fri)	Frequência Acumulada (Fi)	Frequência Acumulada Relativa (Fri)
0 —10	2	2/9 \cong 22%	2	2/9 \cong 22%
10 —20	4	4/9 \cong 44%	6	6/9 \cong 67%
20 —30	3	3/9 \cong 33%	9	9/9 = 100%
Soma (Σ_i)	9	9/9 = 100%	-	-

Como calcular o número de classes?

Op. 1: regra de Sturges

$$nc = 1 + 3,3 \log n$$

Op. 2: r. quadrada

$$nc = \sqrt{n}$$

Como calcular a amplitude do intervalo de classe(h)?

$$h = \frac{(X_{\max} - X_{\min})}{nc}$$

Formas de apresentação dos dados

Uma variável	Duas ou mais variáveis
- Gráfico de frequência - Gráfico de barras - Histograma - Diagrama de pontos - Polígono de frequência - Curva de frequência - Diagrama de ramos e folhas	- Tabelas - Gráfico de colunas - Gráfico de barras - Gráfico de setores (pizza) - Gráfico de dispersão - Gráfico de linhas - Diagrama de ramos e folhas

Estatística Descritiva

*Quais dados devem ser ordenados? mediana, separatrizes

2. Medidas descritivas

A. Medidas de Dispersão:

A1 Absoluta: amplitude total, amplitude interquartilica, desvio quartil, desvio médio, desvio padrão.

A2 Relativa: coeficiente de variação, coeficiente de variação quartil.

B. Medidas de Posição:

B1 Separatrizes: mediana(dados ordenados), decis, quartis, percentis.

B2 Tendência central: média, mediana, moda

C. Medidas de Forma:

C1 Assimetria

C2 Curtose

Transformação uniforme de dados:

Medidas de posição: acompanham a transformação do conjunto de dados (+, -, *, /)

Média, Mediana, Moda, Separatrizes (quartis e decis)

Medidas de dispersão: + e - não afetam

- variância e desvio padrão: * e / sofre modificações (pendente)

Coeficiente de variação: sofre alteração + ou -, * e / não afetam
coeficiente de variação

A. Medidas de dispersão

A1 Absolutas: amplitude total, amplitude interquartilica, desvio quartil, desvio padrão, desvio médio.

Coeficiente de variação interquartilica	Amplitude/ intervalo interquartil (Aq): Amplitude semi-quartilica (As):	Desvio Quartil (Dq)
$\frac{Q3-Q1}{Q3+Q1} = \frac{Aq}{Q3+Q1}$	$A_q = Q_3 - Q_1$ $A_s = \text{mediana}$	$D_q = (Q_3 - Q_1) / 2$
Amplitude total (At):	Desvio:	Desvio médio
$A_t = X_{\max} - X_{\min}$	Desvio = $x_i - \mu$	$D_m = \sum x_i - \mu / n$
Variação:	Desvio padrão: amostra	Desvio padrão população
$\sigma^2 = \sum (x_i - \bar{x})^2 / n$ ou $s^2 = \sum (x_i - \bar{x})^2 / n - 1$	$D_p = \sqrt{\sum (x_i - \bar{x})^2 / n - 1}$	$D_p = \sqrt{\sum (x_i - \bar{x})^2 / n}$

Pq o desvio médio deve ser apresentado em módulo?

x _i	$X_i - \mu$	Desvio	Desvio
2	$2 - 6,4 = -4,4$	-4,4	Soma -6,2
5	$5 - 6,4 = -1,4$	-1,4	
6	$6 - 6,4 = -0,4$	-0,4	
9	$9 - 6,4 = 2,6$	2,6	Soma +6,2
10	$10 - 6,4 = 3,6$	3,6	
Soma ($X_i - \mu$)	0	0	

Ao somarmos os valores de desvio o resultado será 0. Por conta disso, calcularemos o valor do desvio com módulo para dar continuidade ao cálculo.

x _i	Desvio	$ X_i - \mu $
2	-4,4	4,4
5	-1,4	1,4
6	-0,4	0,4
9	2,6	2,6
10	3,6	3,6
Somatório	0	12,4

$$D_m = 12,4 / 5 = 2,48$$

B2 Relativas: coeficiente de variação, coeficiente de variação quartil.

Coeficiente de variação: analisa a dispersão mas não depende de

unidade de medida

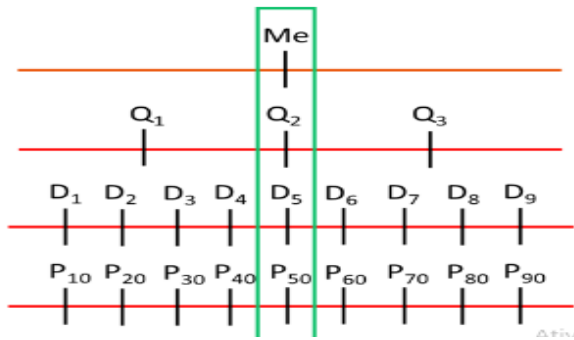
$$Cv = \sigma/\mu$$

Coefficiente de variação interquartil (pendente)

B. Medidas de posição

B.1 Separatrizes

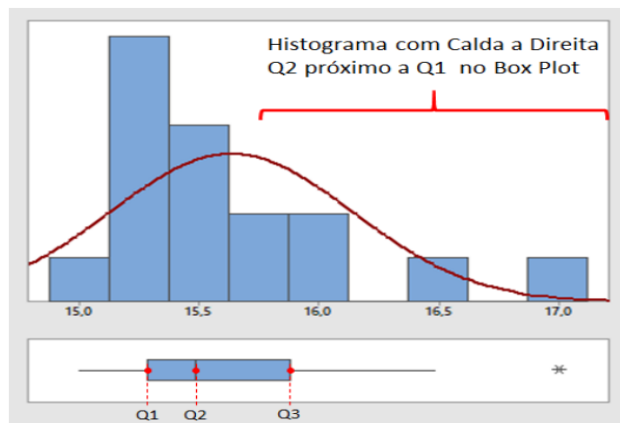
> Quartis (Q)	> Decis (D)	> Percentis (P)
Q1: $0.25(n+1) = (n+1)/4$	D1: $0.10(n+1)$	P5: $0.05(n+1)$
Q2: $0.50(n+1) = (n+1)/2$	D2: $0.20(n+1)$	P20: $0.20(n+1)$
Q3: $0.75(n+1) = 3(n+1)/4$	D3: $0.30(n+1)$	P50: $0.50(n+1)$



Comparativo entre as separatrizes

$$Me = Q2 = D5 = P50$$

Figura 2. Relação de Histograma com Box Plot



B. Medidas de posição

B.2. Tendência Central

a. Média

Considere $x = \{1, 2, 3\}$

Aritmética

$$\bar{x} = (1+2+3) / 3 = 6/3 = 2$$

Geométrica:

$$\bar{G} = \sqrt[3]{1 \cdot 2 \cdot 3} = \sqrt[3]{6} \approx 1,817$$

Harmônica:

$$\bar{H} = \frac{27}{\frac{1}{1} + \frac{1}{2} + \frac{1}{3}} = \frac{27}{1,5} = 18$$

Importante $\bar{x} \geq \bar{G} \geq \bar{H}$

Média aritmética para dados brutos

Cálculo padrão

Média aritmética para dados ponderados

Cálculo	Nota	Peso
$\bar{X}_p = \frac{7 \cdot 3 + 6 \cdot 3 + 8 \cdot 2 + 9 \cdot 1 + 7 \cdot 1}{3+3+2+1+1}$	7	3
	6	3
	8	2
	9	1
	7	1

Média aritmética para dados agrupados

Passo 1: calcular o ponto médio para cada classe	Passo 2: repetir cálculo aplicado para dados ponderados
$Pm_1 = \frac{10+0}{2} = 5$	$\bar{X} = \frac{5 \cdot 2 + 15 \cdot 4 + 25 \cdot 3}{9}$
$Pm_2 = \frac{20+10}{2} = 15$	$\bar{X} = 16,11 \text{ kg/semana}$
$Pm_3 = \frac{30+20}{2} = 25$	

V. observado	Freq Absoluta	Freq Relativa
0-10	2	2/9
10-20	4	4/9
20-30	3	3/9
Somatório	9	9/9

b. Mediana: valor central

Dados devem estar ordenados

Mediana para dados brutos

$$X = \{15, 20, 10, 30, 20, 15, 0, 5, 15\} \quad n = 9$$

$$X = \{0, 5, 10, 15, 15, 15, 20, 20, 30\} \quad (\text{Rol crescente})$$

$$Me = 15$$

Mediana para dados ponderados

Identificar a classe central através da fórmula $n/2$. O valor da moda é o valor observado para a classe.

V. Observado	F. Acumulada	F. Relativa
0	1	11%
5	2	22%
10	3	33%
15 (Mediana)	6	67%
20	8	89%
30	9	100%

Medida para dados agrupados

Passo 1: Identificar a classe central através da fórmula $n/2$. O valor da moda é o valor observado para a classe.

V. Observado	F. Acumulada	F. Relativa
0 - 10	2	22%
10 - 20 (Classe modal)	6	67%
20 - 30	9	100%

Passo 2: calcular o valor da interpolação linear para calcular a mediana.

$$\frac{20-10}{6-2} = \frac{Me-10}{4,5-2}$$

$$\frac{10}{4} = \frac{Me-10}{2,5}$$

$$2,5 = \frac{Me-10}{2,5}$$

$$2,5 \cdot 2,5 = Me - 10$$

$$6,25 + 10 = Me$$

$$Me = 16,25$$

c. Moda: valor que mais se repete

Valor que mais se repete.

>Unimodal:

$$x = \{2, 3, 4, 4, 4, 5, 8\}$$

>Bimodal:

$$x = \{2, 3, 4, 4, 4, 6, 7, 7, 7\}$$

>Amodal:

$$x = \{2, 4, 7, 8, 9, 10, 15\} \quad ()$$

Moda para dados brutos

Passo 1: observar o valor que mais se repete no conjunto de dados

$$x = \{0, 5, 10, 15, 15, 15, 20, 20, 30\}$$

$$Mo = 15 \text{ kg/semana}$$

Moda para dados ponderados

Observar o valor que mais se repete no conjunto de dados através da frequência absoluta (f_i). A moda é o valor observado relativo à frequência

Moda para dados agrupados

Passo 1: observar o valor que mais se repete no conjunto de dados através da frequência absoluta (f_i).

Passo 2: calcula-se o valor pontual da moda através. Há quatro métodos possíveis:

$$Mo = \frac{20+30}{2} = 25 \text{ kg/semana}$$

Moda de Czuber

$$Mo = Li + h \frac{f_{modas} - f_{ant}}{2f_{modas} - (f_{ant} + f_{post})}$$

Moda de Pearson

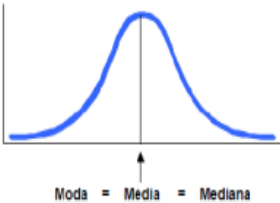
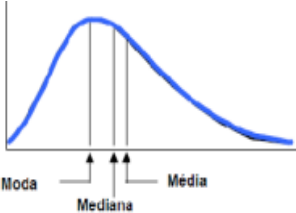
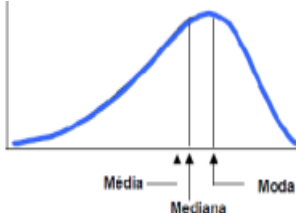
$$Mo = 3Me - 2\bar{x}$$

Moda de King

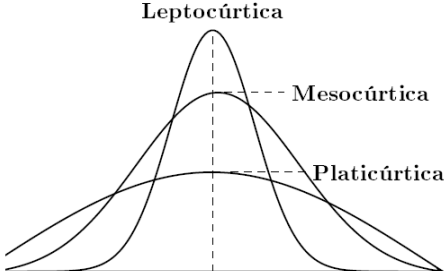
$$Mo = Li + h \frac{f_{post}}{(f_{ant} + f_{post})}$$

C. Medidas de forma

C.1 Assimetria (As)

Distribuição simétrica 	Distribuição Assimétrica a/ direita (positiva) 	Distribuição Assimétrica a/ esquerda (negativa) 	Assimetria $\frac{Q1 + Q3 - 2Me}{Q3 - Q1}$ As > 0 (ass. positiva) As = 0 (simétrica) As < 0 (ass. negativa)
$Q2 - Q1 = Q3 - Q2$	$Q2 - Q1 < Q3 - Q2$	$Q3 - Q2 > Q2 - Q1$	

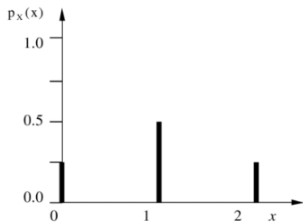
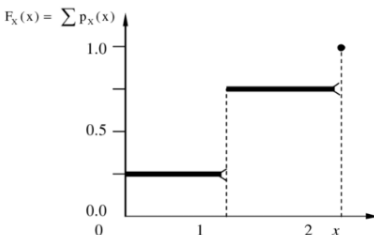
C.2 Curtose

Mesocúrtica	Cálculo de Curtose (C)
	$C = \frac{Q3 - Q1}{2(P90 - P10)}$ ou $C = \frac{Q3 - Q1}{2(D9 - D1)}$

Distribuição de probabilidade discreta (adicionar geométrica)

Distribuição	Caso geral	Bernoulli	Binomial (*)	Poisson	Hipergeométrica (*)
Como identificar	Situações que podem ser quantificadas de forma discreta, isto é, valores inteiros.	1 tentativa 2 resultados (0 / 1)	"n" Bernoulli -independentes	"n" Bernoulli -independentes -intervalo contínuo (tempo/espaco)	"n" Bernoulli -dependentes
Valor Esperado (média)	$E(x) = \sum_{i=1}^n x_i \cdot P(x_i)$	$E(x) = p$	$E(x) = np$	$E(x) = \lambda$	$E(x) = np$ $E(x) = sp$
Variância	$E(x) = E(x_i - E(x))^2 \cdot P(x_i)$ $E(x) = E(x^2) - [E(x)]^2$	$E(x) = pq$	$E(x) = npq$	$Var(x) = E(x) = \lambda$	$Var(x) = npq \left(\frac{N-n}{N-1}\right)$ + fator correção
Desvio padrão	$\sqrt[3]{\text{variância}}$	$\sqrt[3]{\text{variância}}$	$\sqrt[3]{\text{variância}}$	$\sqrt[3]{\text{variância}}$	$\sqrt[3]{\text{variância}}$
Coef. de Variação	$\frac{\text{Desvio padrão}}{E(x)} \cdot 100$	$\frac{\text{Desvio padrão}}{E(x)} \cdot 100$	$\frac{\text{Desvio padrão}}{E(x)} \cdot 100$	$\frac{\text{Desvio padrão}}{E(x)} \cdot 100$	$\frac{\text{Desvio padrão}}{E(x)} \cdot 100$
Função de Probabilidade	$\sum_{i=1}^n P(x = x_i)$	$p^s \cdot q^{1-s}$	$C_{n,s} p^s \cdot (q)^{n-k}$	$e^{-\lambda} \lambda^k$ ----- k! e = 2,72 λ = frequência média fornecida k = número de ocorrências em um intervalo variável (valor buscado)	$\frac{C_x' \cdot C_{n-x}^{N-r}}{C_n^N}$

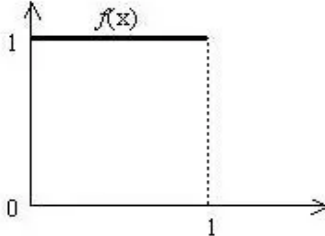
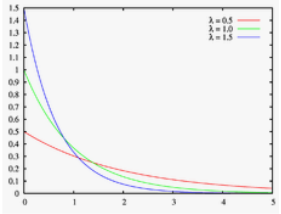
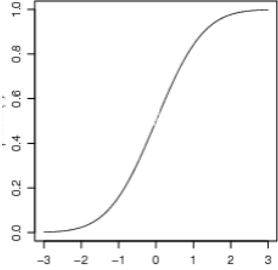
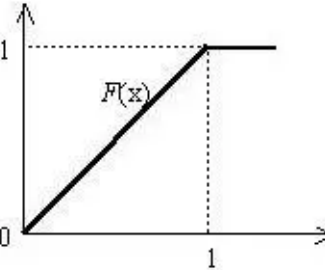
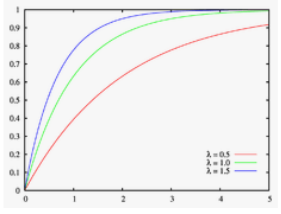
Para todos os casos

Função Massa de probabilidade	Função Distribuição Acumulada de Probabilidade
	

Distribuição de probabilidade contínua

Distribuição	Caso geral	uniforme	Normal	Exponencial
Como identificar	Situações que podem ser quantificadas de forma contínua, isto é, valores decimais.	Probabilidade igualmente distribuída no intervalo		<ul style="list-style-type: none"> - Similar: a Poisson, diferem porque queremos saber o 'tempo/ unidade contínua' em vez de qtd. de ocorrências - Eventos independentes - número de ocorrências por intervalo de tempo deve ser cte λ tempo médio entre ocorrências
Valor Esperado (média)		$\frac{X_{max} - X_{min}}{2}$	$E(x) = \mu$	$E(x) = Dp(x) = 1/\lambda$ - $E(x) = dp(x)$
Variância		$\frac{(X_{max} - X_{min})^2}{12}$	$E(x) = \sigma^2$	$E(x) = 1/\lambda^2$
Desvio padrão	$\sqrt[2]{\text{variância}}$	$\sqrt[2]{\text{variância}}$	$\sqrt[2]{\text{variância}}$	$\sqrt[2]{\text{variância}} = 1/\lambda$
Coefficiente de variação	$\frac{\text{Desvio padrão}}{E(x)} \cdot 100$	$\frac{\text{Desvio padrão}}{E(x)} \cdot 100$	$\frac{\text{Desvio padrão}}{E(x)} \cdot 100$	$\frac{\text{Desvio padrão}}{E(x)} \cdot 100 = \frac{1/\lambda}{1/\lambda} = 1$
Função de Probabilidade (densidade)	$f(x) = P(x) / x_{max} - x_{min}$ Probabilidade no ponto = 0 Caso haja necessidade de calcular a probabilidade de intervalo utiliza-se integral e derivada. Normalmente, utilizamos valores tabelados.	$\frac{1}{X_{max} - X_{min}}$	$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$	$f(x) = \lambda e^{-\lambda x}$ Função acumulada de probabilidade: $f(x) = 1 - e^{-\lambda x}$

PARA DISTRIBUIÇÃO UNIFORME E NORMAL

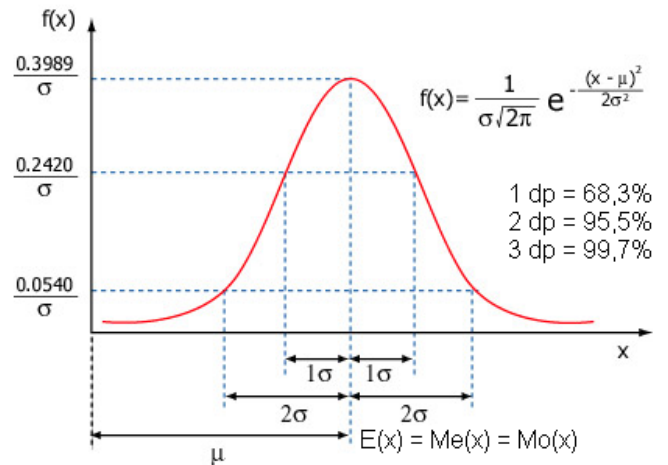
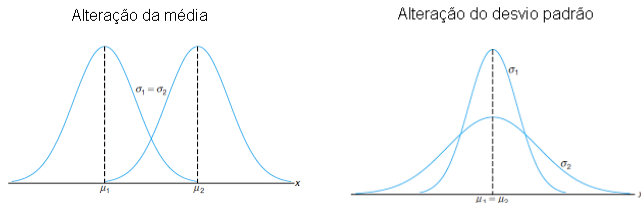
	UNIFORME	EXPONENCIAL
Função Densidade de Probabilidade Atenção: em variáveis contínuas não há Função Massa de probabilidade	Função Densidade de Probabilidade 	Função Densidade de Probabilidade 
Função de distribuição acumulada de probabilidade 	Função densidade acumulada de probabilidade 	Função densidade acumulada de probabilidade 

Atenção: além das distribuições apresentadas acima, há outras distribuições como T de Student, Qui Quadrado, entre outros.

Sobre a distribuição normal

Além de ser utilizada para explicar a probabilidade de variáveis contínuas, também explica intervalos de segurança, teste de hipótese e tamanho amostral.

- Forma de sino
- Simétrica ($\bar{X} = Me = Mo$)
- $P(x < \mu) = P(x > \mu) = 50\%$
- É sempre unimodal
- Mesocúrtica
- Quartis equidistantes ($Q2 - Q1 = Q3 - Q2$)
- É especificada pela média e desvio padrão (raiz quadrada da variância)



Como aplicar a transformação normal nos dados:

1 - Transformar os valores para distribuição de probabilidade normal padrão (Z)

Z: indica o número de desvios padrões a partir da média

$$Z = \frac{X - \mu}{\sigma}$$

Quando aplicamos o teste Z
 $\mu = 0$ e $\sigma = 1$

2 - Identificar a área de Z sobre a distribuição normal.

3 - Encontrar valor tabelado que representa a probabilidade da área determinada.

O valor final de Z indica a quantidade de desvios padrão em relação à média.

Testes para verificar a normalidade dos dados:

Númericos:

-Shapiro-Wilk (limite de 5.000 amostras)

-Kolmogorov-Smirnov (Teste de lilliefors - independe do tamanho de amostras)

-Anderson-Darling

(14:31 -

[https://www.udemy.com/course/estatistica-para-analise-de-dados-com](https://www.udemy.com/course/estatistica-para-analise-de-dados-com-python/learn/lecture/25824976?start=90#questions)

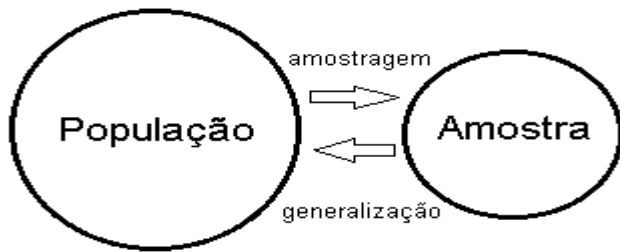
[python/learn/lecture/25824976?start=90#questions](https://www.udemy.com/course/estatistica-para-analise-de-dados-com-python/learn/lecture/25824976?start=90#questions))

Gráficos:

-Histograma

-QQplot

Estatística Inferencial



proporção amostral = probabilidade amostral

Técnicas de amostragem

- Não probabilística (não aleatória)

- Probabilística (aleatória)

a. Simples

- Todas as unidades amostrais devem ter a mesma probabilidade de serem sorteadas

- Seleção das amostras deve permitir ser realizada com ou sem reposição

b. Estratificada

- proporcional

- não proporcional

* Apenas deve-se usar o processo de amostragem estratificada quando houver diferença significativa entre as médias dos estratos, caso contrário, utiliza-se amostragem simples. Nesse caso a amostragem estratificada aumentará a precisão em relação a amostragem simples.

c. Conglomerados: selecionar 2 de 10 turmas de uma escola

d. Sistemática: seleciona elementos com base em frequência pré-definida

$$\text{Fração amostral} = \frac{\text{tamanho da amostra}}{\text{tamanho da população}}$$

Estadística inferencial

OBS: Necessário saber qual o tipo de distribuição que os dados seguem

a - Estudar distribuição amostral

b - Calcular estimadores

c - Calcular intervalo de confiança

d - Calcular tamanho da amostra

e - Calcular erro amostral

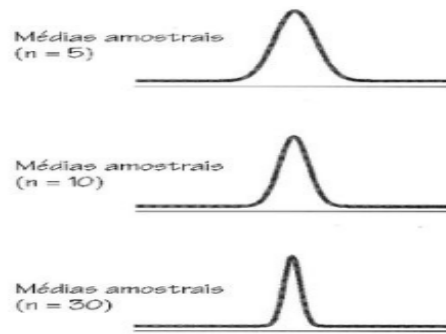
f - Teste de hipótese

Importante: como trabalhamos com amostras, os valores de \bar{x} , s^2 , \hat{p} , \hat{s} são variáveis aleatórias

Importante: as variáveis aleatórias são quantitativas, não aceitam dados descritivos.

Teorema do Limite central

A amostra de dados tende a distribuição normal, mesmo que a população apresente outra distribuição



∴ Não é necessário conhecer a distribuição da população para realizar inferências sobre a amostra.

	Parâmetro população		Estimador amostra
1. Média	méd. população	=	méd. amostra
2. Variância da média (rever e comparar com variância da população)	var. população	≠	var. amostra
Variância da média das amostras	$\sigma^2 = \sum (x_i - \bar{x})/n$ (var. população)	≠	var. amostra = σ^2/n (variância da média da amostra)
	$\sigma^2 = \sum (x_i - \bar{x})/n$	≠	$s_x^2 = \sigma^2/n$ (variância da média da amostra)
3. Desvio padrão da média	$\sigma = \sqrt{\sum (x_i - \bar{x})/n}$	≠	$s = \sigma/\sqrt{n}$ (desvio padrão da média das amostras) "erro média amostral"
4. Probabilidade e proporção "estimativa pontual"	p = interesse/total p = probabilidade	=	p = interesse/total p = proporção
4.1 Esperança	esperança população	=	esperança amostra
Proporção			
4.2 Variância	var = (p.q)	≠	var = (p.q)/n "variância da proporção amostral"
4.3 Desvio padrão	dp = $\sqrt{p.q}$	≠	dp = $\sqrt{\frac{p.q}{n}} = \frac{\sqrt{p.q}}{\sqrt{n}}$ "erro padrão da proporção amostral" "erro da 'estimativa de probabilidade!'"
Para o valor máximo da proporção amostral adota-se p = 50%			

TEOREMA DO LIMITE CENTRAL X LEI DOS GRANDES NÚMEROS

Teorema do limite central: quanto maior a amostra, mais os dados se

aproximam da distribuição normal

>> Converge em distribuição

Lei dos grandes números: quanto maior a amostra, maior a probabilidade de se acertar o parâmetro populacional

>> Forte: converge quase certamente para a média

>> Fraca: converge em probabilidade para a média

∴ A lei dos grandes números não tem relação com a distribuição normal

Converter média amostral para forma normalizada

Dist. Normal	Dist. T de Student
$Z = (X - \mu_x) / \sigma_x$ $Z = (X - \mu_x) / \sigma / \sqrt{n}$	$T = (X - \mu_x) / s_x$ $T = (X - \mu_x) / s / \sqrt{n}$ <p>obs: $s = n - 1$</p>
<p>Quando usar a normal</p> <ul style="list-style-type: none"> - Desvio padrão populacional é conhecido - desvio padrão populacional é desconhecido, porém amostra > 30 	<p>Quando usar t de Student:</p> <ul style="list-style-type: none"> - desvio padrão desconhecido - amostras de $n < 30$ <p>! Difere da normal porque possui distribuição de probabilidade para cada grau de liberdade (GL)</p>
<p>Quando usar T de Student?</p> <ul style="list-style-type: none"> - dp populacional desconhecido - amostra de tamanho pequeno <p>IMPORTANTE: tanto a distribuição normal como a distribuição t de Student tem média = 0. Na dist. normal o desvio padrão é igual a 1.</p>	

$$E = Z \frac{\sigma}{\sqrt{n}}$$

$$\sqrt{n} \cdot E = Z \sigma$$

$$\sqrt{n} \cdot E = Z \sigma$$

$$\sqrt{n} \cdot E = \frac{Z \cdot \sigma}{E}$$

$$n = \left(\frac{Z \cdot \sigma}{E} \right)^2$$

$$n = Z^2 \cdot \frac{p(1-p)}{e^2}$$

Intervalo de confiança



P/ a média amostral (z t)		P/ a proporção amostral (z)	
$\bar{X} \pm Z \frac{\sigma}{\sqrt{n}}$		$\hat{p} \pm Z \sqrt{\frac{\hat{p} \cdot (1-\hat{p})}{n}}$	
P/média amostral	Intervalo de confiança p/ média	Erro de estimativa (margem de erro) p/média (*)	Erro padrão p/ média (*)
- P/ qlqr tamanho de amostra - Dp. pop. conhecido	$\bar{X} \pm Z \frac{\sigma}{\sqrt{n}}$	$Z \frac{\sigma}{\sqrt{n}}$	$\frac{\sigma}{\sqrt{n}}$
- amostra >= 30 - Dp. pop. desc.	$\bar{X} \pm Z \frac{s}{\sqrt{n}}$	$Z \frac{s}{\sqrt{n}}$	$\frac{s}{\sqrt{n}}$
- amostra < 30 - Dp. pop. desc.	$\bar{X} \pm t \frac{s}{\sqrt{n}}$	$t \frac{s}{\sqrt{n}}$	$\frac{s}{\sqrt{n}}$
P/ proporção amostral (b)	Intervalo de confiança p/ média	Erro de estimativa (margem de erro) p/média (*)	Erro padrão (*)
	$\hat{p} \pm Z \sqrt{\frac{\hat{p} \cdot (1-\hat{p})}{n}}$	$Z \sqrt{\frac{\hat{p} \cdot (1-\hat{p})}{n}}$	$\sqrt{\frac{\hat{p} \cdot (1-\hat{p})}{n}}$

Tamanho da amostra

Teste de Hipótese para uma amostra (utiliza valor Z ou T)

- As hipóteses são sempre calculadas para parâmetros populacionais

- Não precisa ser realizado com toda a população
- Similar a intervalo de confiança, no entanto, retorna sim ou não
- Coleta amostra e verifica se há compatibilidade com o valor informado

H_0 = Hipótese nula = deve conter sinal de igualdade obrigatoriamente

$H_0 = \mu = M_0$

H_1 = Hipótese alternativa = desigualdade

$\mu \neq M_0$ - bilateral

$\mu > M_0$ - Unilateral à direita

$\mu < M_0$ - unilateral à esquerda

- Nível de significância: probabilidade do erro. Definido pelo investigador, assim sendo, trata-se de parâmetro subjetivo. Área determinada pelo Z calculado. (chamado nível crítico). "Probabilidade máxima permitida para cometer erro do tipo I"

- Nível de confiança: complemento do nível de significância (lembrando que o total é 1)

Teste de hipótese

Para a média
(permite Z ou T)

$$Z_{cal} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

Z_{cal} = Estatística do teste

Para a proporção
(permite apenas Z)

$$t_{cal} = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

$$Z_{cal} = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$

Etapas para a aplicação:

1. Calcula **Z tabelado**: define a área de rejeição/aceitação (costuma-se usar o valor de 5%, Z= 1,96)
2. Calcula **Z calculado** para a média: verifica posicionamento do valor dentro da distribuição normal

Tipos de erro no teste de hipótese:

Erro do tipo 1: rejeita H_0 quando deveria aceitar

$[\alpha = 1 - \text{nível de confiança}]$

Erro do tipo 2: aceita H_0 quando deveria rejeitar

$[\beta = 1 - \text{potência do teste}]$

Valor - P (Φ)

Importante: o valor Φ também chamado de: probabilidade de significância do teste ou nível descritivo do teste (Varia de 0 a 1)
O nível de significância é o α (determinado pelo pesquisador)

Nível crítico: apresenta nível de aceitação/ rejeição.

Valor-p: é a área projetada pelo Z calculado. Representa a probabilidade de a amostra pertencer a distribuição do valor populacional testado. Pode ser utilizado como uma forma alternativa para que não se calcule a estatística de teste.

Quanto menor o valor de P, maior a probabilidade de rejeitar H_0

área valor p > área nível de significância: aceita H_0

área valor p \leq área nível de significância: não aceita H_0 .

Para encontrar o valor P busca-se na tabela t/z

Teste de hipótese Qui - Quadrado

Variáveis qualitativas ou discretas

- Observações independentes entre si
- Observação contabilizada por contagem, frequência e proporções
- Cada observação pertence a apenas uma classe
- Não pode ser aplicada a amostra inferior a 5 observações por classe

Hipóteses do teste Qui - Quadrado

H_0 : $F_{obs} = F_{esp}$

H_1 : $F_{obs} \neq F_{esp}$

χ tabelado = encontrado na tabel

χ calculado = calculado pela fórmula

Teste qui-quadrado: não utiliza Z e nem t. Utiliza Qui-Quadrado

$$\chi^2 = \sum Z_i^2$$

- Qui-quadrado tabelado

- Qui-quadrado calculado

F_o : frequência observada

F_e : frequência esperada

H_0 : $F_o = F_e$

H_1 : $F_o \neq F_e$

Tipos de teste Qui-quadrado

Aplicado apenas para variável categórica ou discreta, não pode ser aplicado em amostras muito pequenas. $N < 5$

Teste de adequação de ajustamento: 1 variável qualitativa

Teste de independência: 2 ou + variáveis qualitativas

$$\chi^2_{tabelado} = GL$$

$$GL = n-1$$

$$\chi^2_{tabelado} = GL$$

$$GL = (col-1)(lin-1)$$

$$\chi^2_{calculado} = \sum \frac{(F_o - F_e)^2}{F_e}$$

Desvio = $F_o - F_e$ (no entanto, o resultado dessa equação seria zero. Então, elevamos ao quadrado.

Desvio quadrado = $(F_o - F_e)^2$

Proporção: $(F_o - F_e)^2 / F_e$

1 - seleciona a variável com menor número de classes, calcula o % e desenvolve a tabela de valor esperado.

Variáveis bidimensionais

Covariância: + infinito a - infinito
 Correlação: +1 a -1 (indica grau de similaridade entre variáveis diferentes)
 Variância: apenas assume valores positivos
 Coeficiente de determinação R²: apresenta apenas variação de 0 a 1

Distribuição de probabilidade conjunta e marginal

x	0	1	Total
1	1/31/4	2/31/4	1/4
2	1/31/4	2/31/4	1/4
3	1/4
4	1/4
Total	1/3	2/3	

	Conjunta
	Marginal

Probabilidade de variáveis discretas independentes

Como saber se as variáveis são independentes?

A probabilidade conjunta deve ser o resultado da multiplicação das respectivas marginais.

Nesse caso $P(x) = P(x) \cdot p(y)$

Probabilidade de variáveis discretas	
dependentes	independentes
$P_{(x=3 y=1)} = \frac{P(x=3 \text{ e } y=1)}{P(y=1)}$	$P_{(x,y)} = p(x) \cdot p(y)$ p/ todos os pares x e y. Em variáveis independentes a covariância é igual a zero.
Como saber se as variáveis são dependentes? A probabilidade conjunta não é o resultado da multiplicação das respectivas marginais.	

Covariância [valor positivo ou negativo]

-: direção contrária

0: variáveis independentes OU desvios se anulam

+: mesma direção

* Não possui valor mínimo nem máximo

Medida descritiva que indica a relação de dependência entre duas variáveis. Pode ser aplicada tanto em estatística descritiva quanto inferencial.

Covariância como variável descritiva: a covariância é similar a variância, a diferença é que aqui analisamos 2 ou mais variáveis, enquanto na variância analisamos apenas uma.

x	y	$(x - \bar{x})$	$(y - \bar{y})$	$(x - \bar{x}) \cdot (y - \bar{y})$
2	5	-3	-4	12

3	7	-2	-2	4
6	8	1	-1	-1
9	16	4	7	28
				43

$$\text{Cov}(x, y) = \frac{43}{4} \text{ unidade: dominador/ denominador}$$

$$\text{Cov}(x, y) = \frac{\sum(x - \bar{x}) \cdot (y - \bar{y})}{n}$$

Lembre-se de verificar se é n ou se é n-1

Fórmula alternativa da covariância

$\text{Cov}(X, Y) = \text{Média Produto } XY - \text{Produto da média } X \text{ e } Y$

x	y	Produto X e Y
2	5	10
3	7	21
6	8	48
9	16	144

OBS: se fosse amostra seria n-1

Para corrigir a unidade de medida usamos a correlação.

Lembre-se: se as variáveis são independentes se a covariância é 0!

Covariância = 0 não indica necessariamente que as variáveis são independentes!

Qual o problema da covariância? variáveis estão em unidades diferentes. Para solucionar o problema calculamos o coeficiente de correlação.

Covariância como variável aleatória

Correlação (r) [valor entre -1 e +1]: indica a "força" que mantêm unidas duas variáveis.

$$r = \frac{\text{Cov}(X, Y)}{s_x \cdot s_y}$$

$$r = \frac{\sum(x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \cdot \sum(y_i - \bar{y})^2}}$$

Coeficiente de correlação (r)	Correlação Positiva	Coeficiente de correlação (r)	Correlação Negativa
$r = 1$	Perfeita	$r = -1$	Perfeita
$0,95 \leq r < 1$	Muito forte	$-0,95 \leq r < -1$	Muito forte
$0,8 \leq r < 0,95$	Forte	$-0,8 \leq r < -0,95$	Forte
$0,5 \leq r < 0,8$	Moderada	$-0,5 \leq r < -0,8$	Moderada
$0 \leq r < 0,5$	Fraca	$0 \leq r < -0,5$	Fraca

1: variação da variável A é explicada por B

0: não há variação entre as variáveis

-1: correlação negativa entre as variáveis

"Força que mantêm unidos dois conjuntos de valores, por ser adimensional podemos comparar os valores (o que não ocorre na covariância)"

0,0019 a 0,19 - muito fraca
 0,20 a 0,39 - fraca
 0,40 a 0,69 - moderada
 0,70 a 0,89 - forte
 0,90 a 1 - muito forte

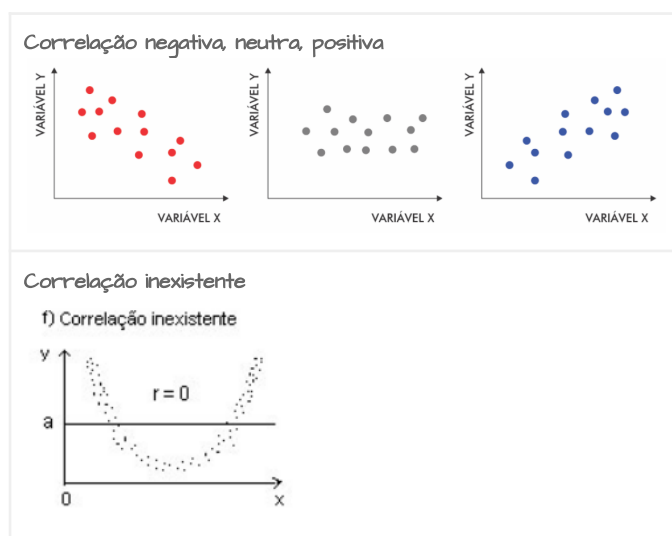
IMPORTANTE: correlação não é causa !

$$r(x, y) = \frac{\text{Cov}(x, y)}{s_x s_y} \text{ [remove unidades de medidas]}$$

O: não possuem dependência ou a correlação não é linear
 s: desvio padrão de x e desvio padrão de y

A covariância não permite realizar comparações por conta da unidade de medida, nesse caso, utilizamos a correlação para poder comparar as variáveis. Dessa forma, como a correlação é adimensional, podemos realizar comparações.

-Costuma usar gráfico de dispersão/ correlação



correlação positiva

Pesquisar sobre root mean square error (RMSE)

- Interpretação facilitada

Pendente:

coeficiente de determinação: $(\text{coeficiente de correlação})^2$

Correlações paramétricas:

Pearson

Correlação não paramétrica

Sperman: permite uso para dados lineares e não lineares

Kendall: utilizado para amostras de até 30 elementos ou para população com grande quantidade de empates

Teste de hipótese para correlação linear

Teste t

$$t = \frac{r}{\sigma_r} = \frac{r}{\sqrt{\frac{1-r^2}{n-2}}}$$

Graus de liberdade

$$gl = n - 2$$

Coefficiente de Determinação:

Porcentagem da variação de y que pode ser explicada pela relação entre x e y

$$r^2 = \frac{\text{Variação encontrada}}{\text{Variação total}}$$

$$r^2 = \rho^2$$

$$y = m \cdot x + b$$

Coefficientes

$$m = \frac{\sum(x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sum(y_i - \bar{y})^2}$$

$$b = \bar{y} - m\bar{x}$$

Cálculo do Coeficiente de Spearman

$$r_R = 1 - \frac{6\sum d_i^2}{n(n^2-1)}$$

n = número amostras.

di = diferença de alcance de cada elemento.

Coeficiente de correlação (r_R)	Correlação Positiva	Coeficiente de correlação (r_R)	Correlação Negativa
$r_R = 1$	Perfeita	$r_R = -1$	Perfeita
$0,95 \leq r_R < 1$	Muito forte	$-0,95 \leq r_R < -1$	Muito forte
$0,8 \leq r_R < 0,95$	Forte	$-0,8 \leq r_R < -0,95$	Forte
$0,5 \leq r_R < 0,8$	Moderada	$-0,5 \leq r_R < -0,8$	Moderada
$0 \leq r_R < 0,5$	Fraca	$0 \leq r_R < -0,5$	Fraca

Cálculo do Coeficiente de Kendall

$$\tau = \frac{(x_i > x_j \text{ e } y_i > y_j \text{ ou se } x_i < x_j \text{ e } y_i < y_j) - (\text{quantidade de pares discordantes})}{n(n-1)/2}$$

Coeficiente de correlação (τ)	Correlação Positiva	Coeficiente de correlação (τ)	Correlação Negativa
$\tau = 1$	Perfeita	$\tau = -1$	Perfeita
$0,95 \leq \tau < 1$	Muito forte	$-0,95 \leq \tau < -1$	Muito forte
$0,8 \leq \tau < 0,95$	Forte	$-0,8 \leq \tau < -0,95$	Forte
$0,5 \leq \tau < 0,8$	Moderada	$-0,5 \leq \tau < -0,8$	Moderada
$0 \leq \tau < 0,5$	Fraca	$0 \leq \tau < -0,5$	Fraca

Análise de Regressão Linear

(explica o formato da regressão)

REGRESSÃO LINEAR SIMPLES

Variável X: 1 independente
Variável Y: 1 dependente
 $y = a + bx + e$

REGRESSÃO LINEAR MÚLTIPLA

Variável X: n variáveis x
Variável Y: 1 variável y
 $y = a + b_1 + b_2 + \dots + b_n + e$

OBS: ao comparar mais de 3 variáveis, não é possível representar graficamente

a: constante da regressão, coeficiente linear, intercepto
valor de x que intercepta y

b: coeficiente de regressão, coeficiente angular, inclinação da reta. Indica a inclinação da reta representa o crescimento de y a cada uma unidade de x. Indica se a reta será /, \ ou -.

Como calcular o valor de A e B?

Estimadores dos mínimos quadrados

$$b = \frac{Cov(x,y)}{S_x^2} \text{ unidade x / unidade y}$$

$$a = \bar{Y} - b\bar{X}$$

ERROS DA REGRESSÃO

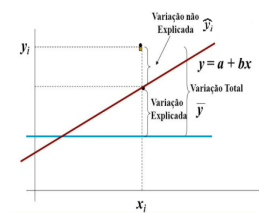
ERROS DA REGRESSÃO

$$e = y_i - \hat{y}_i$$

VARIAÇÃO TOTAL

Explicada: $\hat{y} = a + bx$

Não explicada: $e = y - \hat{y}_i$



Erros da regressão:

$$(y - \bar{y}) = (\hat{y} - \bar{y}) + (y - \hat{y})$$

$$\sum (y - \bar{y})^2 = \sum (\hat{y} - \bar{y})^2 + \sum (y - \hat{y})^2$$

$$SQT = SQE + SQR$$

$$\text{Var. total} = \text{Var. explicada} + \text{Var. erro}$$

Variação = soma dos quadrados

Variação dos erros da regressão

$$se^2 = \frac{\sum (y_i - \hat{y}_i)^2}{n-2}$$

Desvio padrão dos erros da regressão (Erro padrão da estimativa da regressão)

$$se = \sqrt{se^2} = \frac{\sqrt{\sum (y_i - \hat{y}_i)^2}}{\sqrt{n-2}}$$

Teste de hipótese da regressão (usa apenas teste t):

H₀: não existe regressão linear

H₁: existe regressão linear

$$t_{cal} = \frac{\bar{X} - \mu}{s_x} = \frac{b - 0}{s_b} = \frac{b}{s_b} = \frac{se^2}{\sqrt{(x - \bar{x})^2}}$$

Coeficiente de Determinação (r²):

$$r^2 = \frac{SQE}{SQT} = \frac{\text{explicada}}{\text{total}}$$

Análise de variância: (usa f de Snedecor):

$$f \text{ de Snedecor} = \frac{QME(\text{explicada})}{QMR(\text{residual})}$$

Indica qntas vezes a explicada é maior q a residual.

Teste de hipótese dos erros da regressão:

H₀: valor explicado = valor residual ∴ não há linearidade entre x e y

H₁: variância explicada > variância residual

Covariância: unidade medida em min/ kg ∴ não podemos comparar. Para comparar utilizamos a correlação. @

$$\frac{\sum (x - \bar{x}) \cdot (y - \bar{y})}{(n) \text{ ou } (n-1)} = \sum xy - ((\sum x \cdot \sum y)/n)$$

$$\text{Probabilidade} = \sum \epsilon(xy) - \sum \epsilon(x) \cdot \sum \epsilon(y)$$

Correlação: (r): + | a - |

$$r = \frac{Cov(x,y)}{S_x \cdot S_y} =$$

Propriedades da covariância:

(+ e -): não altera

(* e /): * ou /

Variáveis dependentes:

$$\text{Var}(x \pm y) = \text{Var}(x) \pm \text{Var}(y) \pm 2 \text{Cov}(x, y)$$

Variáveis independentes:

$$\text{Var}(x \pm y) = \text{Var}(x) + \text{Var}(y)$$

Propriedades da Correlação (s):

só pode ser alterada pela * ou / de valores negativos.

∴ altera apenas o sinal.

Covariância x Correlação

Teste paramétrico: **analisam variáveis dependentes com distribuição conhecida**. Costuma-se dizer que os testes paramétricos são utilizados para distribuição normal por conta do Teorema do Limite Central.

Pré-requisitos:

a. Dados em intervalo (variável dependente precisa ser numérica, discreta ou contínua)

Dados estão em intervalo		Dados não estão em intervalo	
Cat 1	7,2	Cat 1	Sete
Cat 2	6	Cat 2	Seis
Cat 3	5,4	Cat 3	Cinco

b. Independência dos dados amostrais

c. Normalidade

Testes para verificar normalidade:

>> Kolmogorov-Smirnov (amostra > 50)

>> Shapiro-Wilk (amostra < 50)

>> Anderson Darling

>> entre outros

H₀: distribuição é igual à normal ($p > 0.05$)

H₁: distribuição é diferente da normal ($p < 0.05$)

(em caso de H₁ usa-se testes não paramétricos ou realiza-se a transformação dos dados)

d. Homogeneidade da variância

>> Para verificar homogeneidade aplica-se teste de Levene

H₀: variância é homogênea ($p > 0.05$)

H₁: variância não é homogênea ($p < 0.05$)

(em caso de H₁ pode-se aplicar a correção de Welch)

Teste paramétrico: **analisam variáveis dependentes com distribuição NÃO conhecida**

Grupo Pareado: teste do mesmo grupo em intervalo de tempo ou submetido a diferentes testes. Nesse caso, as amostras devem ser de tamanho igual.

Grupo Independente: grupos iniciais são diferentes, ex: feminino e masculino. Nesse caso as amostras podem ser de tamanhos diferentes.

Testes paramétricos: análise descritiva realizada com média e desvio padrão

Testes não paramétricos: análise descritiva realizada com média (2 quartil), e variabilidade (1 e 3 quartil)

! para calcular o intervalo de confiança é necessário saber a distribuição dos dados.

Variável dependente: resposta

Variável independente: variável de grupamento

! Dados qualitativos não podem ser descritos com média e desvio padrão ∴ só pode ser descrita por mediana, primeiro e terceiro quartil.

! Na análise pareada o grupo deve ser do mesmo tamanho

! Em grupos independentes os grupos não precisam ser do mesmo tamanho

comumente utilizado quando uma suposição do teste paramétrico foi violada. Possuem menor probabilidade para detectar um efeito e não possuem intervalo de confiança. Ideal para dados distorcidos, ou que não estão em distribuição normal ou que apresentam ou

Qual o melhor teste de comparação?	Para uma amostra		
	Testes paramétricos		Testes não paramétricos
	Teste Z Desvio padrão populacional conhecido ou $n \geq 30$	Teste t Desvio padrão populacional desconhecido e $n < 30$	Teste de sinais

Qual o melhor teste de comparação?		VARIÁVEL INDEPENDENTE			
		Qualitativa nominal			
VARIÁVEL QUALITATIVA	grupo	2 grupos		3 grupos ou mais	
	Distribuição	Independentes e não pareados	dependentes /pareados	independente e não pareado	dependentes/ pareados
	Teste	Qui-quadrado/ Fisher	MC Nemar	Qui-quadrado/ Fisher	Q de Cochram
		Qualitativa Ordinal			
DISTRIBUIÇÃO DE PROBABILIDADE	grupo	2 grupos		3 grupos ou mais	
	Distribuição	Independentes e não pareados	dependentes /pareados	Independentes e não pareados	dependentes /pareados
	Teste	<u>Mann - Whitney (versão não paramétrica do T de Student)</u>	<u>Wilcoxon</u>	<u>Kruskal - walls +comparação múltipla</u>	<u>Friedman +comparação múltipla</u>
		Discreta/Continua			
DISTRIBUIÇÃO DE PROBABILIDADE	grupo	2 grupos		3 grupos ou mais	
	Distribuição	Independentes e não pareados	dependentes /pareados	independente e não pareado	dependentes/ pareados
	Teste	- Paramétrico T de Student	-Paramétrico T de Student Pareado	- Paramétrico: ANOVA	- Paramétrico ANOVA de medidas repetidas

ATIVAS	- <u>Não paramétrico</u> <u>Teste de Mann-Whitney ("T de Student não paramétrico")</u>	- <u>Não paramétrico</u> <u>Teste de Wilcoxon ("Teste T pareado não paramétrico")</u>	+comp. múltipla (Tukey) - <u>Não paramétrico</u> <u>Kruska Walls</u> + <u>comp. múltipla (Dunn's)</u>	+comp. múltipla (Tukey) - <u>Não paramétrico</u> <u>Friedman</u> + <u>comp. múltipla (*Rever porque foi usado Tukey em vez de Dunn's - https://www.youtube.com/watch?v=piC-hsYazBK)</u>
	OBS: teste de Mann-Whitney e podem ser aplicados tanto a variáveis quantitativas não paramétricas quanto a variáveis qualitativas ordinais		OBS: teste de Kruska Walls e Friedman podem ser aplicados tanto a variáveis quantitativas não paramétricas quanto a variáveis qualitativas ordinais	

Estudo pendente: Análise de Variância (ANOVA) a dois critérios (SigmaPlot 12.0) +Tukey [Revisar]

Fator 1: 2, 3 ou + grupos; Fator 2: 2, 3 ou + grupos

Alguns testes de comparação múltipla: Tukey, Bonferroni, Dunn, Dunnet, Fisher LSD, Sidak, Student-Newman-Keuls, Duncan, Scheffé

Qual o melhor teste de correlação a ser utilizado?		
Teste de correlação de Pearson	não há variáveis independentes	2 variáveis dependentes quantitativas (apenas 2) *Cada unidade amostral deve ser analisada pelas duas variáveis dependentes
Teste de correlação de Spearman		2 variáveis dependentes, sendo: - quantitativa + qualitativa ordinal ou - qualitativa ordinal + qualitativa ordinal pelo menos uma variável avaliada não tem distribuição normal (ex: qualitativa)

Comparação, correlação ou regressão?		Variável Dependente		
		Qualitativa nominal	Quantitativa nominal	Qualitativa
Var iáv el Ind epe nde nte	Qualitativa nominal	Associação	Comparação não paramétrica	Comparação paramétrica ou Regressão Logística
	Quantitativa ordinal	Regressão logística	Correlação não paramétrica ou Regressão logística ordinal	Comparação não paramétrica ou Regressão simples
	Quantitativa	Regressão logística	Regressão logística ordinal	Correlação paramétrica ou Regressão simples