



# Análise Descritiva

# Descrição x Inferência

- A análise descritiva tem por objetivo sumarizar os dados, tornando a sua visualização mais fácil
- Ela não produz inferências, ou seja, serve apenas para organizar e apresentar informações
  - No entanto, muitas vezes apenas a descrição já pode ser suficiente para resolver um problema de pesquisa

# Descrição x Inferência

- Através da análise descritiva podemos resumir um conjunto grande de informações em poucas linhas
- Na aula de hoje vamos usar a base sen2018.csv
  - Para acessá-la abram o Sigaa
  - Carreguem a base e o pacote "tydeverse"

# Tabelas de Frequência

- Apresenta a frequência contagem de uma variável
  - No tidy, podemos usar a função `count()`
  - Ou a função `table()`, do Rbase

# Tabela de Referência Cruzada (Contingência)

- Tabelas de contingência ou de referência cruzada trazem a frequência (contagem) a partir do cruzamento de duas variáveis
  - Qual o partido que apresentou mais candidaturas de mulheres ao senado em 2018?

# Proporções

- Na tabela anterior, vimos que o PSOL apresentou o maior número de candidaturas de mulheres ao senado
- No entanto, os dados absolutos não são comparáveis entre si, porque cada partido lançou um número distinto de candidaturas
- Podemos tornar os casos comparáveis usando proporções

# Tabelas com 3 variáveis

- Qual o partido que apresentou o maior número de mulheres negras ao senado em 2018?

# Tabelas em HTML

- Uma forma de exportar as tribbles geradas em tidyverse para outros documentos é usando a função `flextable()` do pacote "flextable"
  - Ela gera uma tabela em HTML que pode ser aberta no navegador e copiada para outros documentos (planilhas ou texto)



# Medidas de Tendência Central

- São estatísticas descritivas que representam um valor central em uma distribuição. As medidas mais comuns são a média, a mediana e a moda.
  - Média: retorna o valor médio dos casos  $\rightarrow \bar{X} = \frac{\sum X}{N}$
  - Mediana: retorna a localização da posição média de uma distribuição  $\rightarrow$   
$$\text{Mediana} = \frac{N+1}{2}$$
  - Moda: retorna o valor mais comum em uma distribuição. Não é a frequência e sim o valor!

# Problemas da Média

- A média é uma medida imprecisa quando trabalhamos com dados de cauda pesada
  - Cauda pesada = concentração de casos em um dos lados do histograma
  - Isso acontece quando há muita dispersão nos nossos dados (desigualdade)
- Temperatura média em um deserto ou renda em Belém são exemplos dito.

# Medidas de tendência central

- Vimos que a média, a mediana e a moda da idade coincidem no mesmo valor (55), mas podemos afirmar que esta variável apresenta uma distribuição normal?
- Para fazer esta verificação, vamos empregar o teste de Shapiro-Wilk
  - $H_1$  = a distribuição não é normal
  - $H_0$  = a distribuição é normal

# Medidas de tendência central

Teste de normalidade de Shapiro-Wilk, usado para verificar se os dados seguem uma distribuição normal (gaussiana).

variável que está sendo testada para normalidade

## Shapiro-Wilk normality test

```
data: sen2018$IDADE_DATA_POSSE  
W = 0.98955, p-value = 0.02514
```

Valor da estatística de teste W calculada pelo teste de Shapiro-Wilk para os seus dados. Quanto mais próximo esse valor estiver de 1, mais os dados se assemelham a uma distribuição normal. Neste caso, o valor de W é aproximadamente 0,98955.

Valor-p é igual a 0,02514, o que é menor do que 0,05. Portanto, com um nível de significância de 0,05, você rejeitaria a hipótese nula e concluiria que os dados da variável "IDADE\_DATA\_POSSE" não seguem uma distribuição normal.

# Medidas de Dispersão

- A partir da média, podemos calcular o desvio de cada caso. Esta medida indica a distância de um caso da média, logo é dada por:
  - $x - \bar{x}$
- Para calcular o desvio de uma distribuição, o mais lógico seria somar todos os desvios e dividir pelo N. No entanto  $\sum(x - \bar{x})$  será sempre igual a 0.
- Isso ocorre porque, em uma distribuição, os desvios positivos (acima da média) compensam os desvios negativos (abaixo da média).

# Medidas de Dispersão

- Uma forma de contornar este problema é elevar ao quadrado todos os desvios. Esta medida é a variância ou  $s^2$  ou  $\sigma^2$  (sigma ao quadrado).

- $$s^2 = \frac{\sum (X - \bar{X})^2}{N}$$

$s^2$  representa a variância amostral.

$\sum$  denota a soma.

$X$  representa cada valor individual nos dados.

$\bar{X}$  é a média dos dados.

$N$  é o número de observações na amostra.

A variância é uma medida importante, mas os valores estão em unidades ao quadrado, o que pode ser difícil de interpretar diretamente. Portanto, frequentemente usamos a raiz quadrada da variância, que é chamada de desvio padrão ( $s$  ou  $\sigma$ ), para obter uma medida de dispersão na mesma escala que os dados originais.

# Medidas de Dispersão

- Apesar de ser uma medida de dispersão, a variância nos diz muito pouco sobre a dispersão da distribuição, uma vez que ela não representa o valor exato da média dos desvios. Para isto, simplesmente podemos obter uma raiz quadrada da variância, ou o desvio padrão:

$$s = \sqrt{\frac{\sum (x - \bar{x})^2}{N}}$$

- No R, o desvio padrão pode ser obtido através da função `sd()`

# Medida de Dispersão

- Como interpretar o Desvio Padrão?
  - É uma medida de dispersão, ou seja, quanto menor o seu valor, mais próximo da média se encontram os valores da distribuição.
  - Ele pode indicar a homogeneidade ou heterogeneidade dos dados: valores muito alto indicam que os dados estão muito dispersos, ou seja, variam muito em relação à média.
  - Imagine uma sala medindo 5m x 5m. No meio dela, você coloca uma mesa de centro e passa a distribuir o restante dos móveis a partir dela. A  $\pm 1$  metro da mesa você coloca os sofás, a 2 metros uma estante com livros, a -2m um vaso de flor... O desvio padrão seria o metro, ou a medida que você usa para distribuir os móveis nesta sala. Quanto maior for o desvio padrão, mais longe da mesa de centro os restantes dos móveis estarão posicionados.



# Medidas de dispersão

- O valor da variância é mostrado em forma de notação científica.
  - Podemos alterar isso com a função `options(scipen = 999)`
  - Isso faz com que o R exiba números em seu formato decimal normal, sem usar a notação científica.

# Medidas de dispersão

Exemplo genérico para calcular no quadro:

Todo mês acontece na UFPA a Feira de Agricultura Familiar.

Em uma das barrquinhas tem 3 crianças.

Queremos:

- \*Identificamos as idades das crianças
- \*Saber a média das idades
- \*Identificar a variância e o desvio padrão



$$\text{Média}(\bar{x}) = \frac{\sum \text{valores}}{\text{número de observações}}$$

$$\bar{x} = 12$$

$$\text{Variância}(s^2) = \frac{\sum (\text{valor} - \bar{x})^2}{\text{número de observações}}$$

$$s^2 = \frac{24}{3} = 8$$

$$s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}}$$

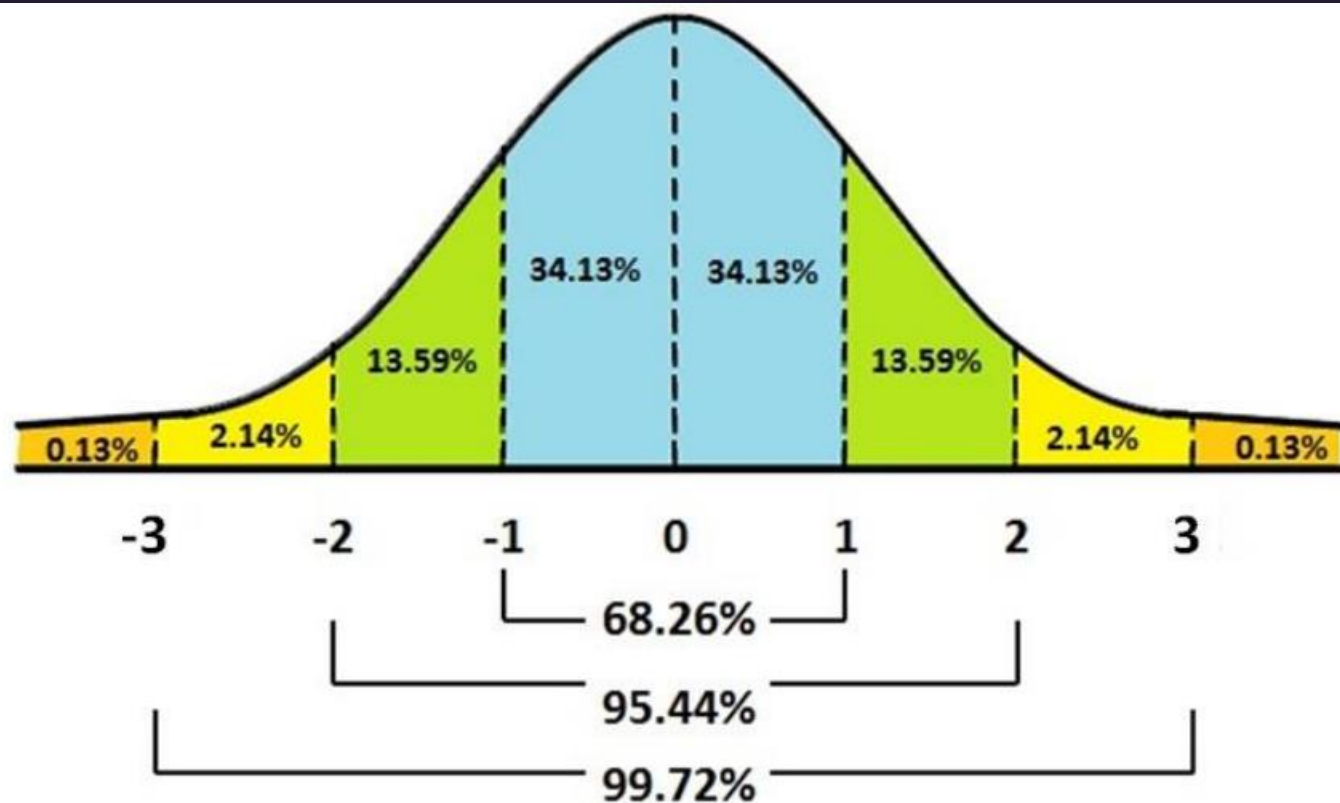
$$s = \sqrt{8} \approx 2,8$$

# Escore Z (também conhecido como escore padrão)

- É uma medida de distância de um valor em relação à média do conjunto
  - Escore aponta a distância em desvios padrão
- Valores de escore Z negativos indicam que o valor produto é inferior à média, já valores positivos significa que estão situados acima da média.

# Escore Z (também conhecido como escore padrão)

- Seus valores oscilam entre  $-3 < Z < +3$  e isto corresponde a 99,72% da área sob a curva da Distribuição Normal.



# Escore Z

Comumente utilizado em estudos de desempenho acadêmico, análises financeiras, controle de qualidade e etc.

- Calcular o escore z é bem simples:

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

- Onde:

- $X$  = valor bruto
- $\bar{X}$  = média da distribuição
- $\sigma$  = desvio padrão

A média das idades é 55 anos

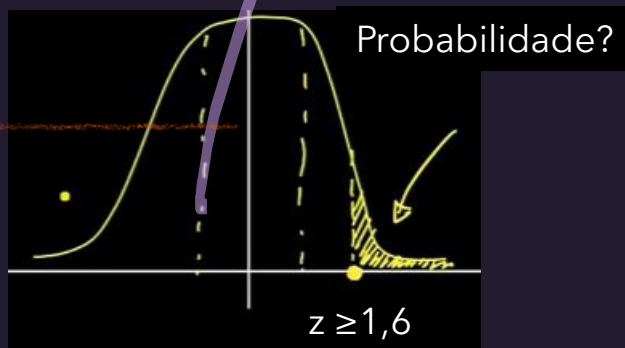
# Escore Z

- Vamos calcular o escore z da idade do senador Jader Barbalho

```
#####Escore Z#####  
#####Escore Z do Jader Barbalho#####  
jader_z <- (74-mean(sen2018$IDADE_DATA_POSSE))/  
  sd(sen2018$IDADE_DATA_POSSE)  
  
jader_z #[1] 1.693606
```

1,6 desvios padrões da média  
 $P(z \geq 1,6)$

\*mostrar o gráfico e ir para a tabela



1. indica que a idade de Jader está cerca de 1.693606 desvios padrão acima da média das idades no conjunto de dados. Ou seja? a idade de Jader é relativamente maior em comparação com a maioria das idades no conjunto de dados.

2. Como o escore Z é positivo, isso indica que a idade de Jader está acima da média das idades no conjunto de dados. Quanto maior o valor do escore Z, maior a diferença entre a idade de Jader e a média.

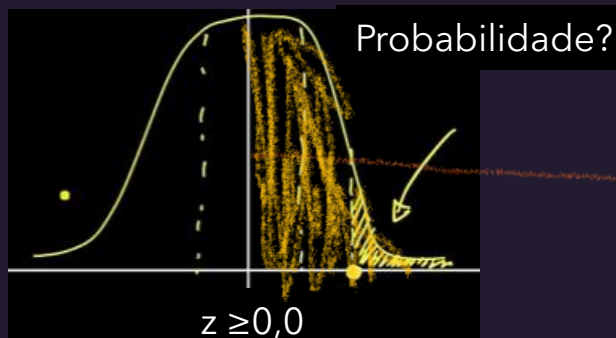
# Escore Z

```
> summary(sen2018$IDADE_DATA_POSSE)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
  27.0   47.0   55.0   55.1   63.0   83.0
> |
```

```
#####Calcular todos os Escores Z da base#####
sen2018$idade_z <- (sen2018$IDADE_DATA_POSSE -
                    mean(sen2018$IDADE_DATA_POSSE))/
                    sd(sen2018$IDADE_DATA_POSSE)

summary(sen2018$idade_z)
mean(sen2018$idade_z)

# > summary(sen2018$idade_z)
# Min.    1st Qu.    Median      Mean    3rd Qu.      Max.
# -2.518651 -0.726202 -0.009222  0.000000  0.707758  2.500208
# > mean(sen2018$idade_z) #Explicar a leitura do valor em notação
# [1] 9.594753e-17
```



Qual a probabilidade abaixo da curva ?

\*tabela







# Risco Relativo

- O Risco Relativo (RR) é uma medida derivada da epidemiologia que ajuda a entender as chances as chances de um evento ocorrer.
- Permite identificar a chance de um evento ocorrer em detrimento de outro (ou todos os outros) eventos
- Através dele podemos identificar a magnitude do fator de risco de um evento.

# Risco Relativo

- O cálculo do risco relativo é feito a partir de tabelas 2x2 (duas linhas e duas colunas) e é dado por:

$$RR = \frac{\textit{Probabilidade de Ocorrência no Fator de Risco}}{\textit{Probabilidade de Ocorrência fora do Fator de Risco}}$$

# Risco Relativo

- Ilustração:

	Ocorrência do Evento	
	SIM	Não
Presença do Fator de Risco	A	B
Ausência do Fator de Risco	C	D

$$RR = \frac{\frac{A}{A+B}}{\frac{C}{C+D}}$$

# Risco Relativo

- Qual o risco relativo de uma mulher ser eleita senadora?

# Risco Relativo

- Como interpretar?
  - Valores negativos indicam uma menor probabilidade de um evento acontecer.
  - No exemplo, as mulheres têm menos 0,69 vezes a chance de serem eleitas, ou
  - $1 - 0,69 = -31\%$  de chances de serem eleitas

# Risco Relativo

- Limite para o Risco Relativo
  - Cervi (2014) aponta que fatores de risco inferiores a  $\pm 50\%$  (risco relativo entre 0,5 e 1,5) são considerados não práticos, ou seja, devem ser desconsiderados.
  - No exemplo anterior, o Risco Relativo é de  $-31\%$ , o que apontaria um fator não relevante.