

**Senac**

**Fecomércio  
Sesc**

## Data Science – Princípios e Técnicas

Outubro  
2024



# Data Science



Onde me encontrar:

<https://www.linkedin.com/in/marco-mialaret-junior/>

e

<https://github.com/MatmJr>

# O conjunto de dados

# Data Science

---

```
import pandas as pd
```

```
url = 'http://personal.tcu.edu/kylewalker/data/colleges.csv'
```

```
data = pd.read_csv(url, encoding = 'latin_1')
```

```
data.shape
```

# Data Science

---

Para responder a essa pergunta, precisamos identificar algumas colunas essenciais e filtrar os dados de acordo. As colunas que manteremos são as seguintes:

- **INSTNM**: Nome da instituição;
- **STABBR**: Estado onde a instituição está localizada;

# Data Science

---

- **PREDDEG:** Tipo principal de diploma concedido pela instituição (códigos: 1 para certificados, 2 para diplomas de associado, 3 para bacharelados e 4 para pós-graduação);
- **CONTROL:** Propriedade da instituição (códigos: 1 para pública sem fins lucrativos, 2 para privada sem fins lucrativos, e 3 para privada com fins lucrativos);

# Data Science

---

- **UGDS:** Número de alunos de graduação matriculados na instituição;
- **UG25abv:** Percentual de alunos de graduação com 25 anos ou mais na instituição.

# Mais um pouco de estatística



# Data Science



---

As medidas de tendência central, por si só, não são suficientes para fornecer uma descrição completa dos dados. É essencial considerar também as medidas de variabilidade, que quantificam a dispersão dos pontos de dados em relação à média.

# Data Science

---

Nesta aula, você aprenderá a identificar e calcular as principais medidas de variabilidade:

- Variância
- Desvio padrão
- Amplitude
- Assimetria
- Curtose

# Data Science



---

## Variância

A **variância** é uma medida estatística que indica a dispersão dos valores em um conjunto de dados. Ela calcula o quão distante cada ponto de dados está da média do conjunto, proporcionando uma visão da variabilidade geral. A variância é geralmente representada pelo símbolo  $s^2$ .

# Data Science

---

## Exemplo 1:

# Data Science



---

## Desvio Padrão

O **desvio padrão** da amostra é uma importante medida de dispersão dos dados, intimamente relacionada à variância da amostra. O desvio padrão, representado por  $s$ , é obtido ao se calcular a raiz quadrada positiva da variância da amostra.

# Data Science

---

Essa medida é frequentemente preferida à variância, pois é expressa na mesma unidade que os dados originais, facilitando a interpretação. Após calcular a variância, o desvio padrão pode ser facilmente determinado usando Python.

# Data Science

---

## Exemplo 2:

# Data Science

---

## Amplitude

Amplitude é uma medida de dispersão estatística que indica a diferença entre o maior e o menor valor em um conjunto de dados. Essa métrica é particularmente útil para entender a escala total dentro da qual os dados variam, oferecendo uma visão rápida da extensão dos valores observados.



# Data Science



---

Para calcular a amplitude, simplesmente subtraímos o valor mínimo do valor máximo encontrado no conjunto de dados. Apesar de sua simplicidade, a amplitude tem limitações significativas, principalmente porque é extremamente sensível a valores atípicos (outliers). Um único valor extremamente alto ou extremamente baixo pode distorcer a percepção da variação geral dos dados.

# Data Science

---

---

## Distribuição de Frequência

A distribuição de frequência é uma ferramenta estatística essencial para analisar e visualizar como os dados estão distribuídos ao longo de diferentes categorias ou intervalos contínuos.

# Data Science

---

Ela permite observar a frequência com que determinados valores ou grupos de valores ocorrem dentro de um conjunto de dados, facilitando a identificação de padrões e a compreensão da estrutura dos dados.

# Data Science



---

A construção de uma distribuição de frequência pode ser realizada através da tabulação dos dados em classes, que podem ser definidas por intervalos de valores. Cada classe contém uma contagem de quantas vezes valores dentro daquele intervalo aparecem no conjunto de dados, conhecida como frequência.

# Data Science



---

A representação visual dessas frequências, geralmente por meio de histogramas (gráficos de barras), proporciona uma compreensão clara e imediata das características principais dos dados, como a concentração de observações em determinados intervalos.

# Data Science

---

**Exemplo 3:** Construa um histograma para visualizar a distribuição de frequência.

Vamos precisar de uma biblioteca nova: matplotlib.

# Data Science

---

```
import matplotlib.pyplot as plt
```

```
plt.hist(variável_numérica)
```

```
plt.title('Título')
```

```
plt.xlabel('Nome do eixo x')
```

```
plt.ylabel('Frequência')
```

```
plt.show()
```



# Data Science



---

## Box-Plot

O box-plot é uma ferramenta gráfica essencial na estatística descritiva, fornecendo uma visão clara da distribuição dos dados, bem como de sua concentração e dispersão. Este gráfico destaca quão distantes os valores extremos estão em relação ao corpo principal dos dados, facilitando a identificação de outliers.

# Data Science



---

Um box-plot é estruturado a partir de cinco valores chave: o mínimo, o primeiro quartil (Q1), a mediana (Q2), o terceiro quartil (Q3) e o máximo. Esses pontos são utilizados para construir o gráfico e comparar a distribuição de diferentes conjuntos de dados.

# Data Science



---

Para desenhar um box-plot, posiciona-se uma caixa retangular sobre um eixo numérico, que pode ser horizontal ou vertical. Os extremos da caixa são definidos pelos primeiro e terceiro quartis, encapsulando aproximadamente 50% dos dados entre eles.

# Data Science

---

O box-plot é uma ferramenta poderosa para uma análise rápida e eficaz, permitindo comparações diretas entre diferentes grupos de dados e ajudando a destacar características como simetria e dispersão de maneira intuitiva e imediata.

# Data Science

---

Vamos precisar de mais uma biblioteca para plotar gráficos, a seaborn

```
import seaborn as sns  
sns.boxplot(x= variável_numérica)
```

# Data Science



---

## Assimetria

A **assimetria** é uma medida que indica o grau de desvio na simetria de uma distribuição de dados. Em termos simples, a assimetria quantifica o quanto a distribuição de uma característica se afasta de uma distribuição perfeitamente simétrica.

# Data Science



**Observação:** Valores negativos de assimetria indicam a presença de uma cauda dominante à esquerda da distribuição, como ilustrado no primeiro conjunto de dados. Por outro lado, valores positivos indicam uma cauda mais longa ou mais espessa à direita, conforme demonstrado no segundo conjunto. Uma assimetria próxima de zero (por exemplo, entre -0,5 e 0,5) sugere que a distribuição é relativamente simétrica.

# Data Science

---

## Exemplo 4:



## Curtose

**Curtose** é uma medida estatística que descreve o "achatamento" ou "pico" de uma distribuição em relação à distribuição normal. Outliers em uma amostra impactam a curtose de forma mais significativa do que afetam a assimetria, pois, numa distribuição simétrica, caudas mais pesadas em ambos os extremos elevam a curtose.

# Data Science



---

Em contraste com a assimetria, onde as caudas opostas podem se neutralizar, na curtose ambas contribuem para seu aumento. Diferentemente da média e do desvio padrão, que são expressos nas mesmas unidades dos dados, e da variância, que é expressa no quadrado dessas unidades, a curtose é uma medida adimensional. Ela representa um coeficiente que indica o grau de achatamento da distribuição dos dados.

# Data Science

---

## Exemplo 5:

# Data Science



---

O padrão de referência para a curtose é a distribuição normal, que tem uma curtose de 3. Por esse motivo, frequentemente utiliza-se o termo excesso de curtose, que é calculado como curtose menos 3.

# Data Science

---

- Mesocúrtica: Uma distribuição é classificada como mesocúrtica quando sua curtose é aproximadamente igual a 3 (ou um excesso de curtose próximo de 0, usando a medida do pandas). Essa distribuição tem uma forma semelhante à distribuição normal em termos de "pico".

# Data Science

---

- Platicúrtica: Uma distribuição com curtose menor que 3 (ou excesso de curtose negativo no pandas) é chamada de platicúrtica. Comparada à distribuição normal, ela tende a ter caudas mais curtas e finas e um pico central mais baixo e largo.

# Data Science

---

- Leptocúrtica: Quando a curtose é maior que 3 (ou excesso de curtose positivo no pandas), a distribuição é classificada como leptocúrtica. Esse tipo de distribuição possui caudas mais longas e gordas, com um pico central mais alto e agudo em comparação com uma distribuição normal.

# Data Science



---

Ao analisar a assimetria e a curtose das idades, observamos que os dados não estão distribuídos de maneira centralizada e apresentam uma forma verticalmente alongada. Para obter uma visualização mais precisa da distribuição dos dados, é útil empregar uma curva conhecida como função densidade de probabilidade.



# Data Science



---

Essa curva suaviza os picos e vales dos histogramas ao passar pelas diferentes classes de dados, oferecendo uma representação contínua e detalhada da distribuição.

# Dúvidas?

---



**Marco Mialaret, MSc**

**Telefone:**

**81 98160 7018**

**E-mail:**

**marcomialaret@gmail.com**

