

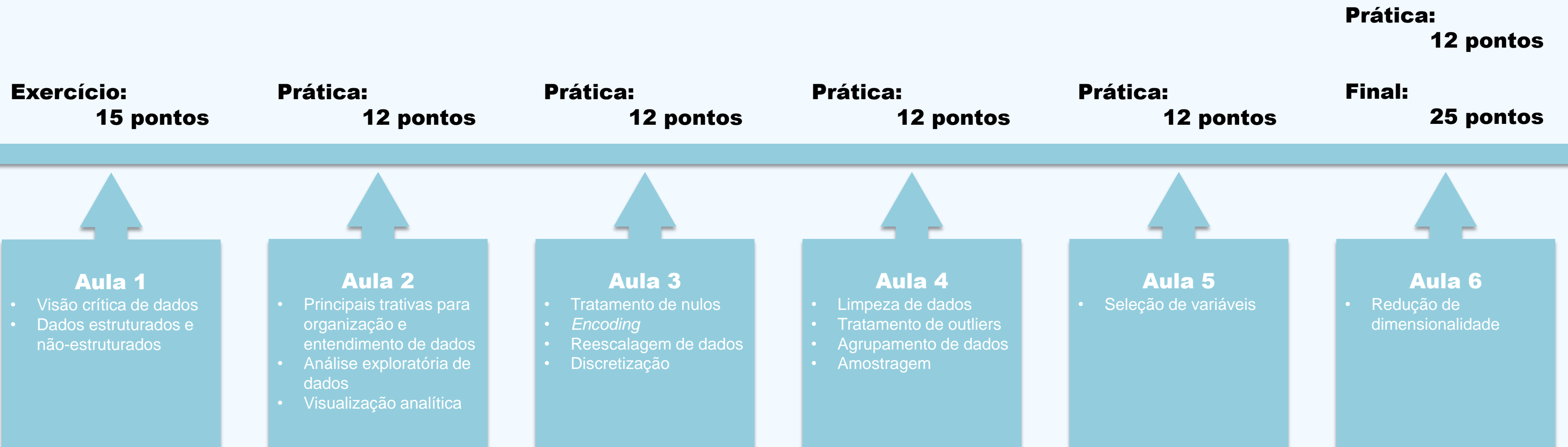
Disciplina:

MODELAGEM E PREPARAÇÃO DE DADOS PARA APRENDIZADO DE MÁQUINA

Professor: Rafael Barroso



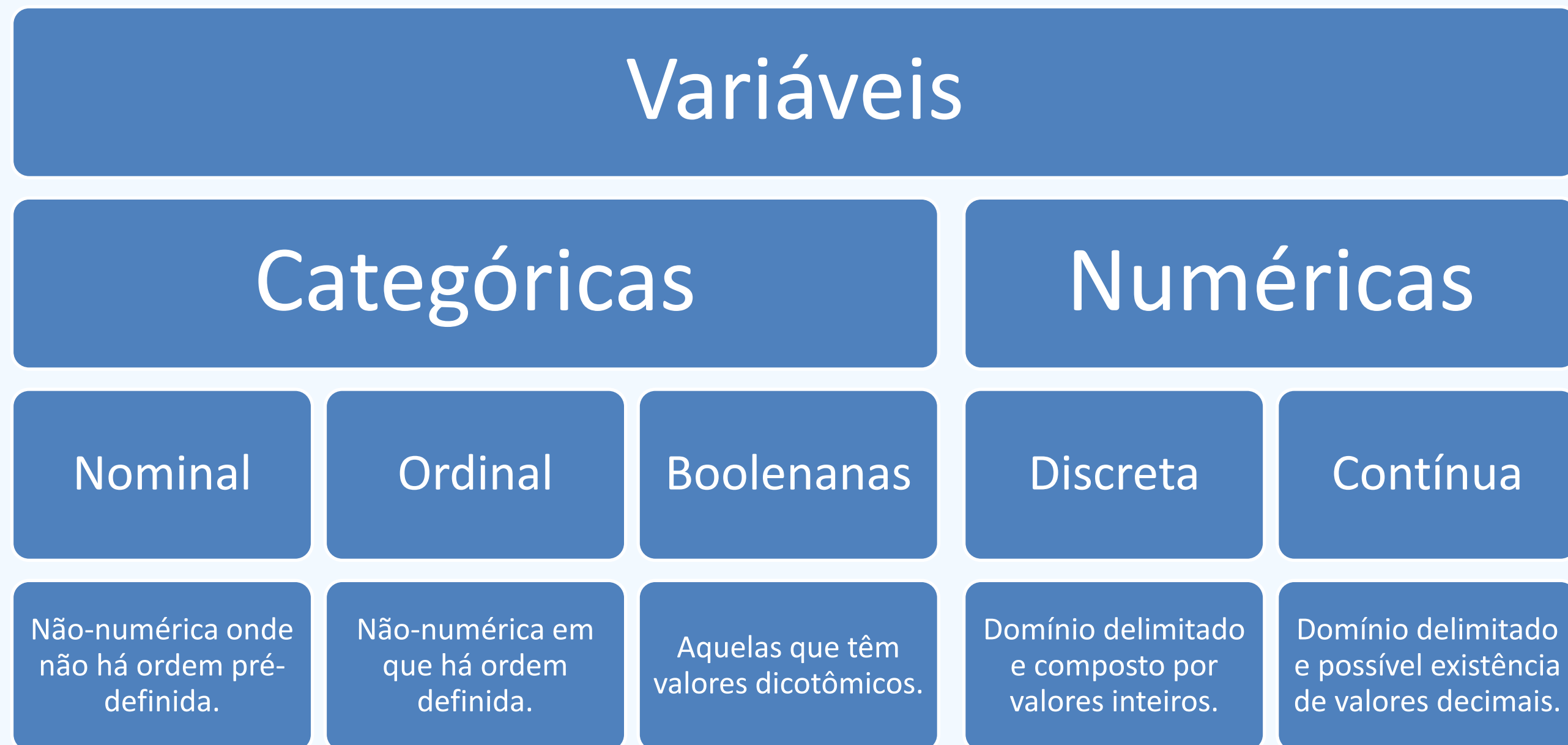
AVALIAÇÃO



PARTE 1

TIPOS DE VARIÁVEIS

TIPOS DE VARIÁVEIS



TIPOS DE VARIÁVEIS

Categórica Nominal

- Cor dos olhos,
- Cor do cabelo,
- Estado de saúde,
- Cidade natal,
- País,
- Estilo musical, etc.

Categórica Ordinal

- Escolaridade,
- Mês (de um dado evento).

Numérica Discreta

- Número de filhos/pessoas,
- Número de prédios,
- Número de assentos, etc.

Numérica Contínua

- Tempo,
- Peso,
- Altura,
- Dinheiro,
- Distância, etc.

PORTE 2

ESTATÍSTICA DESCRITIVA

ESTATÍSTICA DESCRITIVA

MEDIDAS DE TENDÊNCIA

Média

- Valor que resume as observações a um valor central.
- Resume a amostra a um valor que, se repetido n vezes (onde n é o tamanho da amostra), obtém-se o mesmo valor para a soma dos elementos da amostra.
- É fortemente influenciado por valores extremos (mínimos e máximos).
- Pode não ter um significado coerente para o domínio em questão.

Mediana

- Valor que divide as observações da amostra em dois grupos de mesmo tamanho.
- É equivalente ao percentil 50%.
- Demanda que a amostra esteja ordenada para sua identificação.
- Em amostras com número ímpar de elementos, é o próprio elemento central.
- Em amostras com número par de elementos, é a média dos dois elementos centrais.

Moda

- Valor mais recorrente na amostra.
- Pode não existir.

ESTATÍSTICA DESCRITIVA

MEDIDAS DE TENDÊNCIA

12	15	23	32	4	93	14	12	13
----	----	----	----	---	----	----	----	----

Média

Mediana

Moda

$$\bar{x} = \frac{\sum_{i=0}^n x_i}{n}$$

$$\bar{x} = \frac{12+15+23+32+4+93+14+12+13}{9}$$

$$\bar{x} = \frac{218}{9}$$

$$\bar{x} = 24,2$$

4, 12, 12, 13, 14, 15, 23, 32, 93

$$M = 14$$

Número	Frequência
4	1
12	2
13	1
14	1
15	1
23	1
32	1
93	1

ESTATÍSTICA DESCRITIVA

MEDIDAS DE TENDÊNCIA

76	34	64	74	42	23	59	79	26	86
----	----	----	----	----	----	----	----	----	----

Média

Mediana

Moda

ESTATÍSTICA DESCRITIVA

MEDIDAS DE TENDÊNCIA

76	34	64	74	42	23	59	79	26	86
----	----	----	----	----	----	----	----	----	----

Média

Mediana

Moda

$$\bar{x} = \frac{\sum_{i=0}^n x_i}{n}$$

$$\bar{x} = \frac{76+34+64+74+42+23+59+79+26+86}{10}$$

$$\bar{x} = \frac{563}{10}$$

$$\bar{x} = 56,3$$

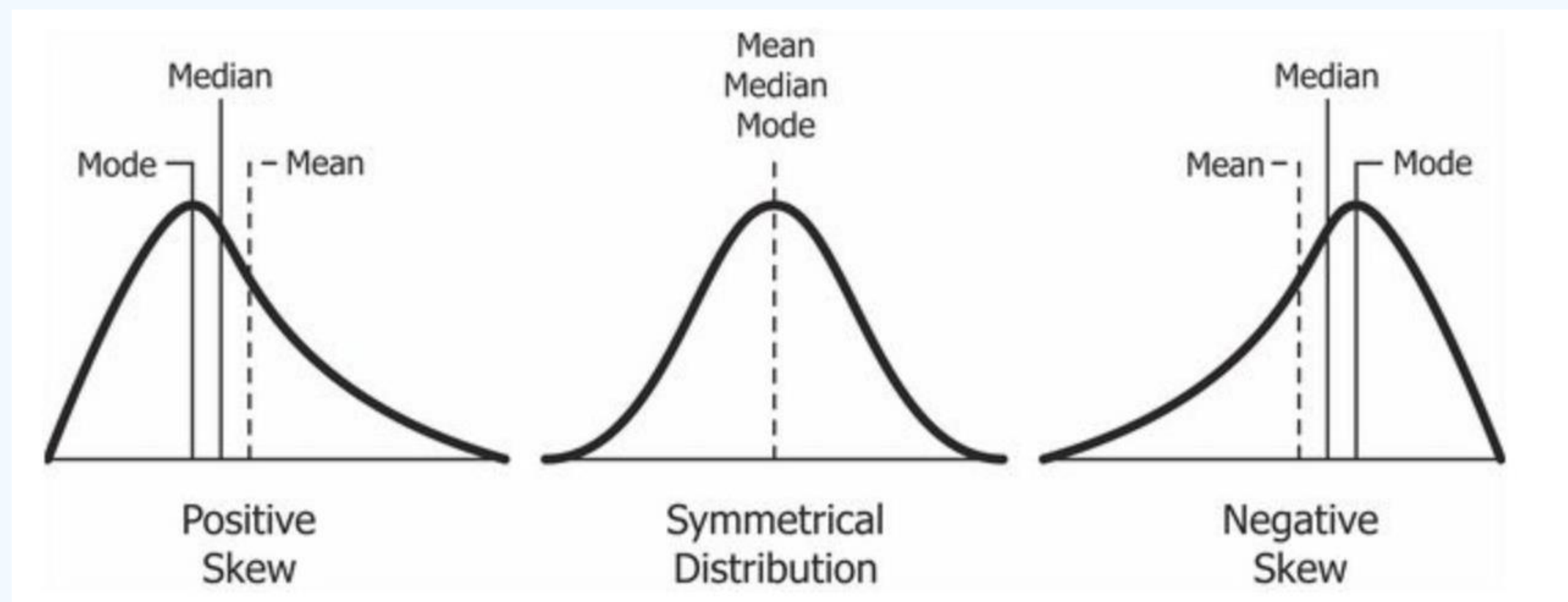
23, 26, 34, 42, 59, 64, 74, 76, 79, 86

$$M = 61,5$$

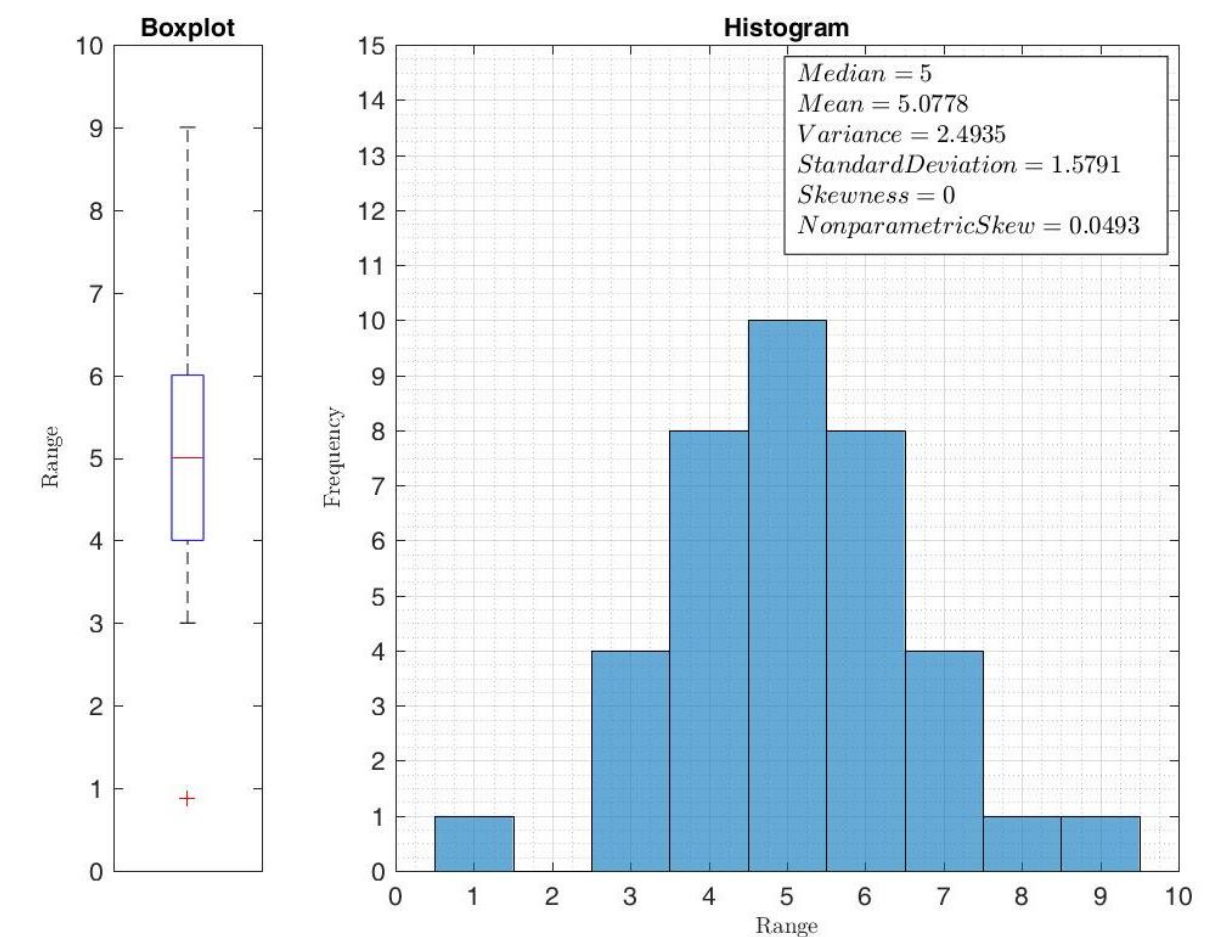
Número	Frequência
76	1
34	1
64	1
74	1
42	1
23	1
59	1
79	1
26	1
86	1

ESTATÍSTICA DESCRITIVA

ASSIMETRIA



Counterexample
Case of an Asymmetric Distribution with Zero Skewness



ESTATÍSTICA DESCRITIVA

MEDIDAS DE VARIABILIDADE

Amplitude

- É a diferença entre os valores máximo e mínimo da amostra de dados.

Desvio Padrão

- Indica o grau de variação no conjunto de dados.
- Valores baixos indicam que os elementos se concentram próximo à média.
- Por conseguinte, valores altos indicam que os dados estão distribuídos/menos concentrados.
- É sempre positivo ou nulo (sempre nulo para constantes).

ESTATÍSTICA DESCRITIVA

MEDIDAS DE TENDÊNCIA

76	34	64	74	42	23	59	79	26	86
----	----	----	----	----	----	----	----	----	----

Amplitude

$$A = x_{\text{máx}} - x_{\text{min}}$$

$$A = 86 - 23$$

$$A = 63$$

Desvio Padrão

$$\sigma = \sqrt{\frac{\sum_{i=0}^n (x_i - \bar{x})^2}{n}}$$

$$\sigma = 22,1$$

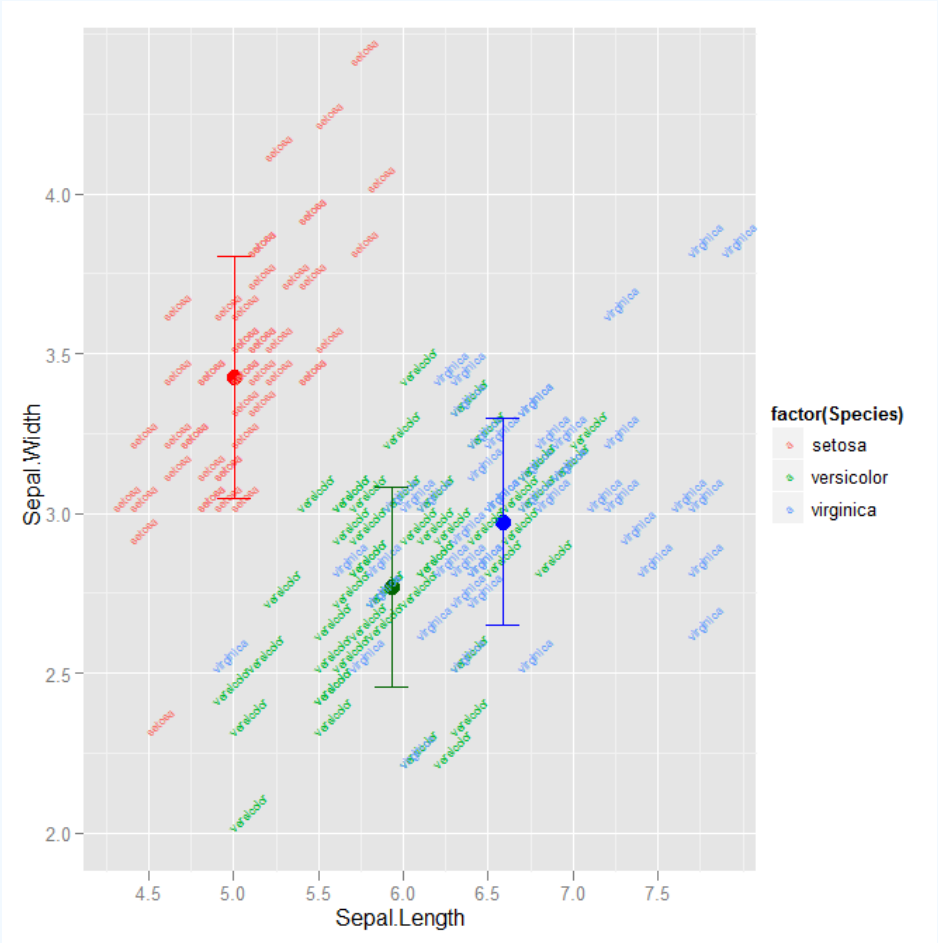
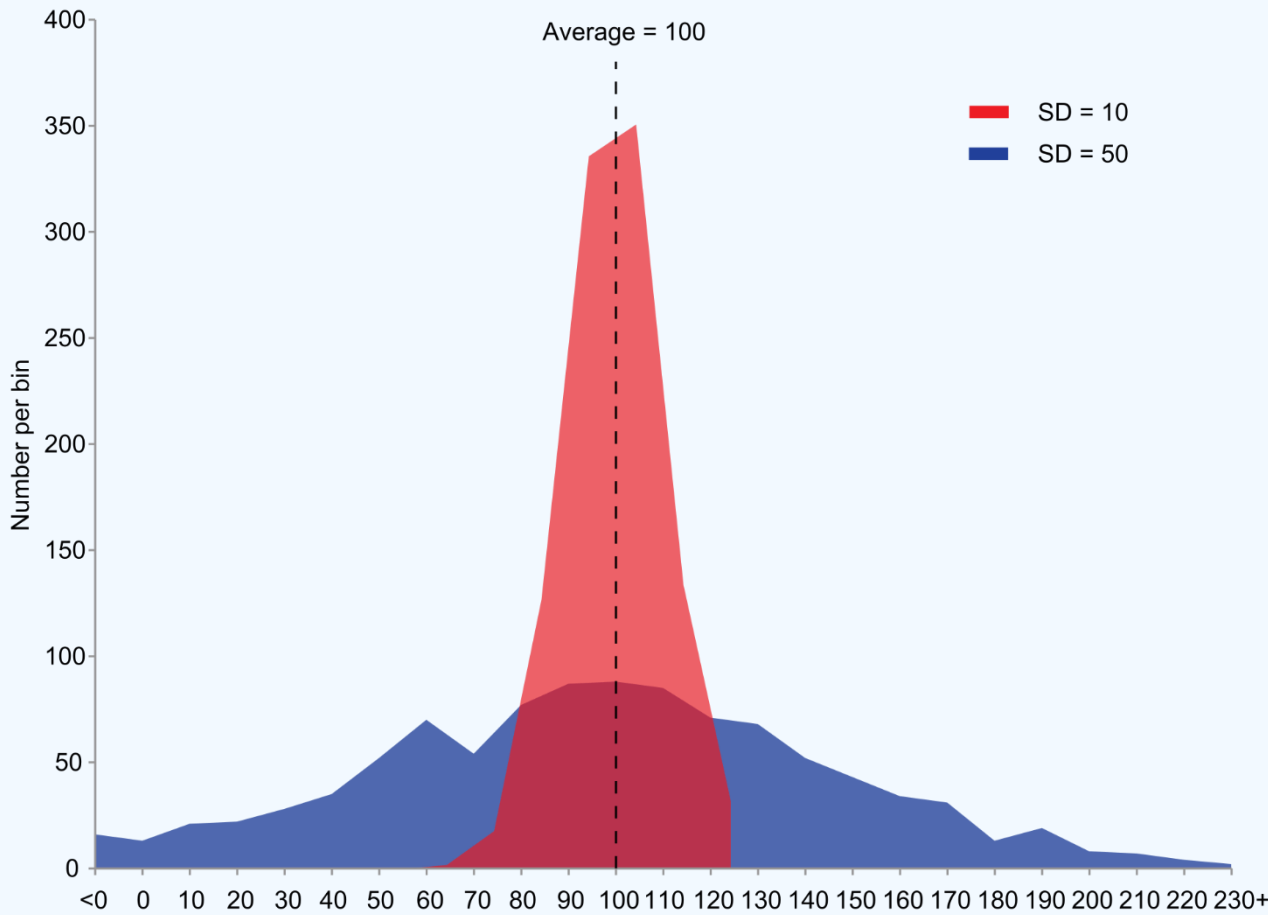
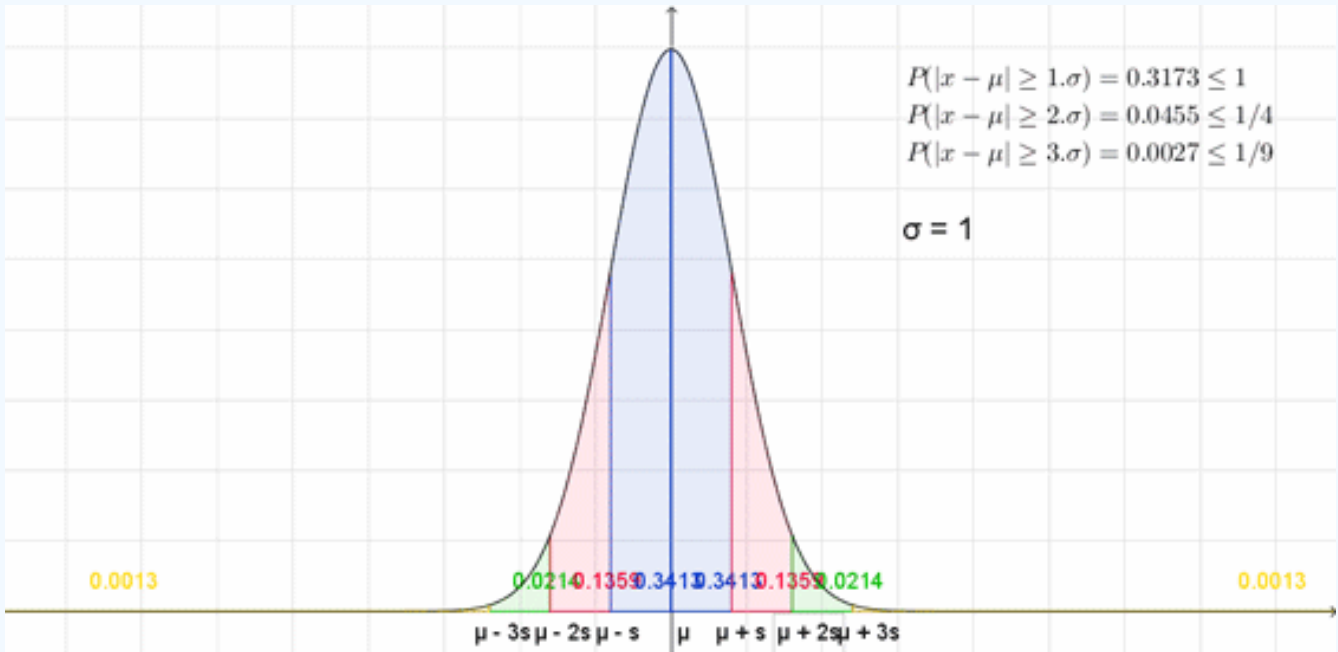
Variância

$$\sigma^2 = \frac{\sum_{i=0}^n (x_i - \bar{x})^2}{n}$$

$$\sigma = 489,4$$

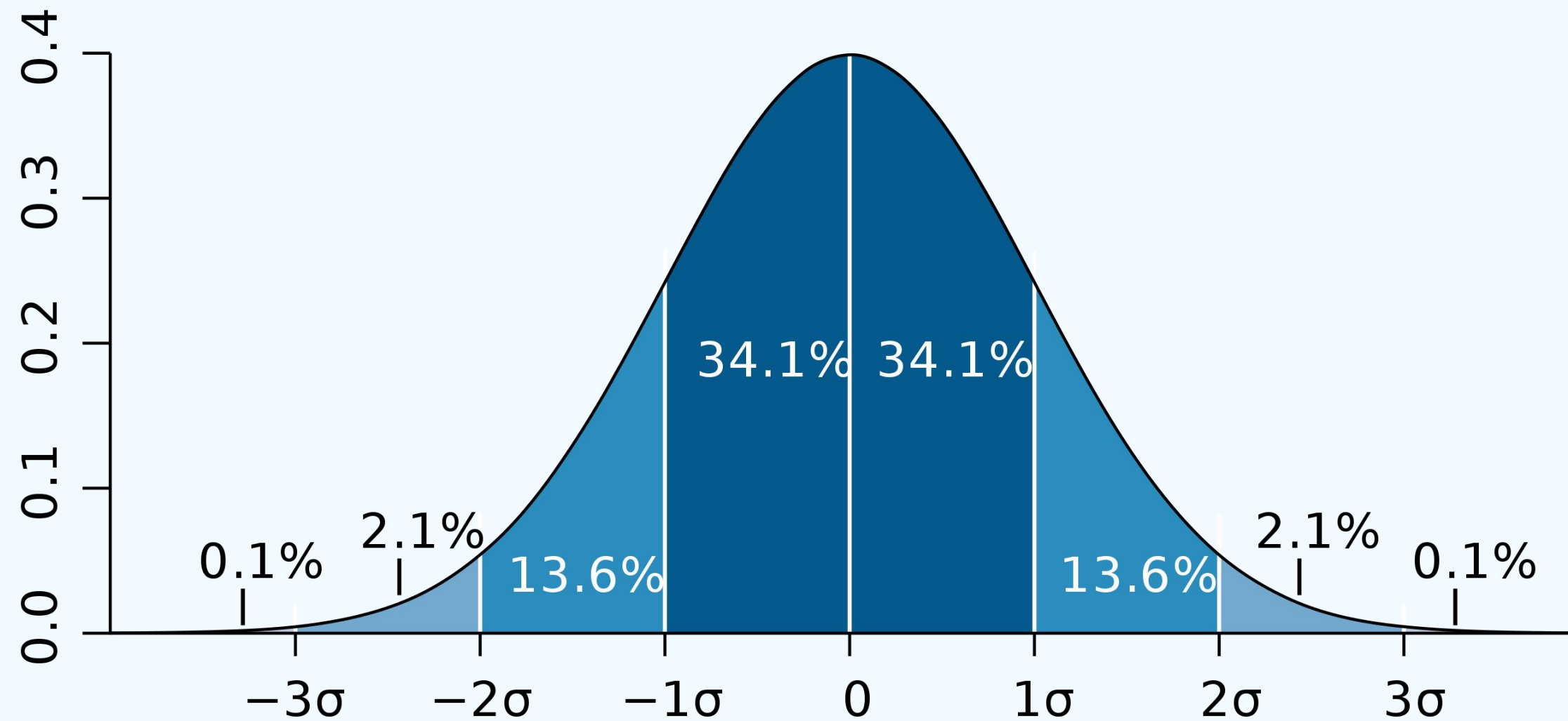
ESTATÍSTICA DESCRITIVA

MEDIDAS DE VARIABILIDADE



ESTATÍSTICA DESCRITIVA

MEDIDAS DE VARIABILIDADE



ESTATÍSTICA DESCRITIVA

MEDIDAS DE POSIÇÃO

- Percentis, quartis, etc.
- Responsáveis por dividir a amostra ordenada de acordo com o percentual da população que desejamos;
- Ex: percentil 10% → equivale ao valor (no domínio

da amostra) que separa os primeiros 10% de observações dos outros 90%;

- Percentis notáveis:
 - 1º quartil → percentil 25%,
 - 2º quartil → percentil 50% → mediana,
 - 3º quartil → percentil 75%.

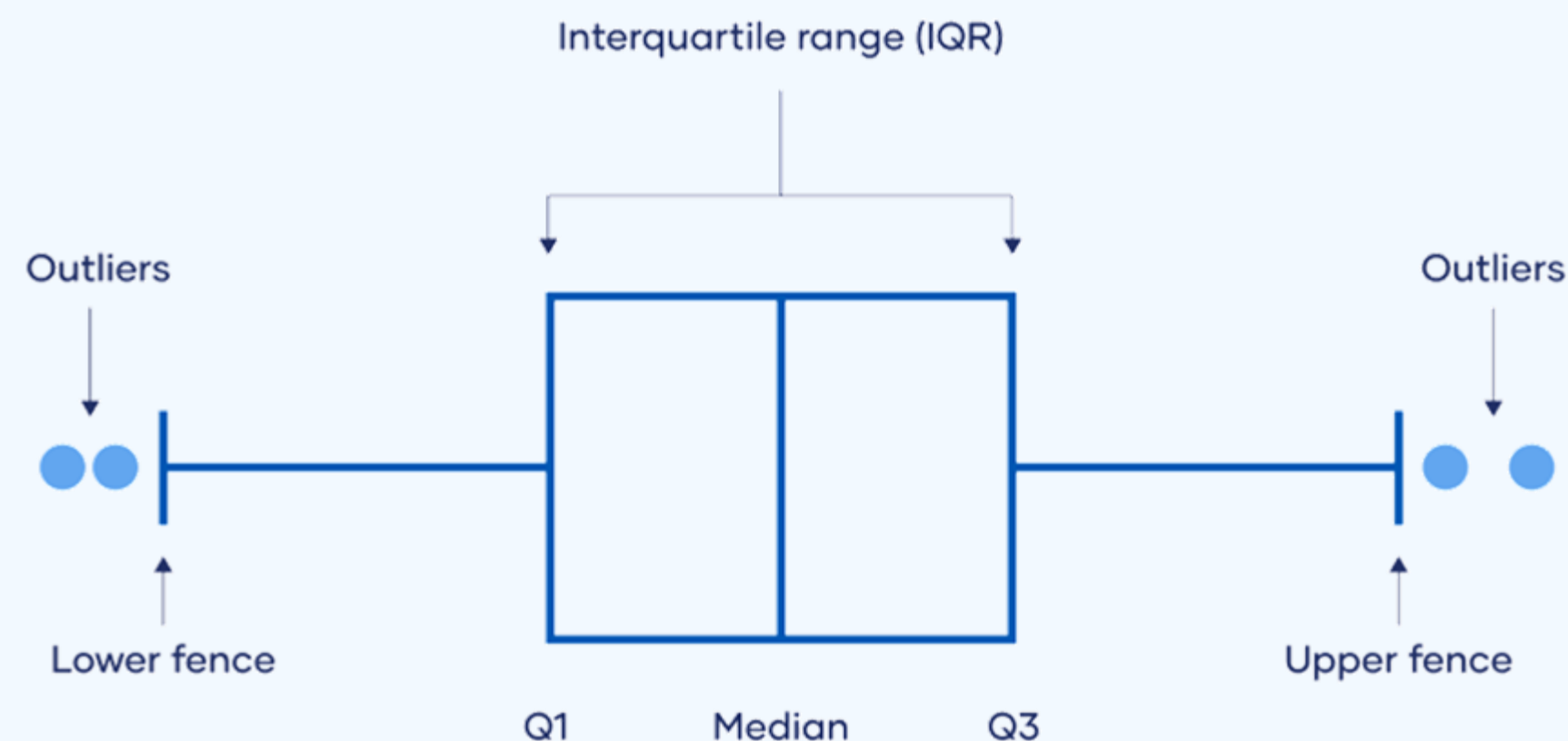
ESTATÍSTICA DESCRITIVA

OUTLIERS

- São valores possíveis, mas incomuns;
- Em alguns casos, podem ser facilmente identificados no conjunto de dados;
- Em outros, pode-se aplicar algum teste estatístico para identifica-los;
- A principal ferramenta se baseia na distância interquartis.

ESTATÍSTICA DESCRITIVA

OUTLIERS



Identificação de Outliers – Distância Interquartis:

$$\text{Limite Inferior} = Q_1 - 1,5 \times IQR$$

$$\text{Limite Superior} = Q_3 + 1,5 \times IQR$$

$$IQR = (Q_3 - Q_1)$$

PARTE 3

ANÁLISE EXPLORATÓRIA

ANÁLISE EXPLORATÓRIA

O QUE É:

- Processo inicial de análise crítica;
- Busca identificar padrões, anomalias, testar hipóteses e validar premissas;
- Tem como objetivo entender os dados antes de começar a usá-los;
- Ajuda a definir como melhor utilizar/trabalhar os dados;
- Além de validar informações conhecidas (ou supostas), pode revelar novos *insights*.

ANÁLISE EXPLORATÓRIA

ABORDAGENS:

