

VOLUME 1:

CIÊNCIA DE DADOS

Fundamentos da Ciência de dados

Eduardo Amadeu Dutra Moresi

Paulo Fernando Marschner



Cátedra UNESCO de
Juventude, Educação
e Sociedade

VOLUME 1:

CIÊNCIA DE DADOS

Fundamentos da Ciência de Dados

Eduardo Amadeu Dutra Moresi

Paulo Fernando Marschner



Brasília, DF

2026

“The designations employed and the presentation of material throughout this publication do not imply the expression of any opinion whatsoever on the part of UNESCO concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. The ideas and opinions expressed in this publication are those of the authors; they are not necessarily those of UNESCO and do not commit the Organization.”



É proibida a reprodução total ou parcial desta publicação, por quaisquer meios, sem autorização prévia, por escrito, da Cátedra Unesco de Juventude, Educação e Sociedade.

A exatidão das informações, conceitos e opiniões é de exclusiva responsabilidade dos autores, os quais também se responsabilizam pelas imagens utilizadas.

Coleção Juventude, Educação e Sociedade

Comitê Editorial:

Geraldo Caliman (Coordenador), Célio da Cunha, Gilvan Charles Cerqueira de Araújo, Jenerton Arlan Schütz, Marta Helena de Freitas, Renato de Oliveira Brito.

Conselho Editorial Consultivo:

Esther Martínéz (Portugal), Azucena Ochoa Cervantes (México), Cristina Costa Lobo (Portugal), Marília Costa Morosini (Brasil), Paulo César Nodari (Brasil).

Capa e diagramação:

Laura Tiemi Fugita.

Equipe técnica TIC em trilhas

Coordenação:

Eduardo Moresi (Coordenador geral), Vilson Hartmann (Coordenador pedagógico), Mário Braga (Coordenador de infraestrutura).

Design Instrucional:

Isolda Gusmão.

M843c Moresi, Eduardo Amadeu Dutra.
Ciência de dados [recurso eletrônico] : fundamentos da ciência de dados :
volume 1 / Eduardo Amadeu Dutra Moresi e Paulo Fernando Marschner. –
Brasília, DF : Universidade Católica de Brasília, 2026.
(Coleção Juventude, Educação e Sociedade).

Inclui referências bibliográficas.
Disponível em: <<https://ucb.catolica.edu.br>>.
ISBN 978-65-87629-79-7

1. Ciência de dados. 2. Estatística aplicada. 3. Aprendizado de máquina. 4.
Inteligência artificial. 5. Tomada de decisão. I. Título. II. Marschner, Paulo
Fernando.

CDU 004.6

Sobre os Autores

Eduardo Amadeu Dutra Moresi

Possui graduação em Engenharia Eletrônica pelo Instituto Militar de Engenharia - IME - Rio de Janeiro (1989), mestrado em Engenharia Elétrica pela Universidade de Brasília - UnB (1994) e doutorado em Ciência da Informação pela UnB (2001). É professor da Universidade Católica de Brasília - UCB, desde 1997. É autor de diversas publicações científicas nacionais e internacionais. Desde 2014 coordena o Programa Apple Developer Academy da UCB. Desde 2023, coordena o Projeto TIC em Trilhas da UCB. Atua como docente e pesquisador nos Programas de Mestrado e Doutorado em Educação e de Mestrado Profissional em Governança, Tecnologia e Inovação.

Pode ser contatado pelos e-mails: moresi@p.ucb.br e eduardo.moresi@gmail.com

ORCID: <https://orcid.org/0000-0001-6058-3883>

LATTES: <http://lattes.cnpq.br/1068744176859901>

Paulo Fernando Marschner

Bacharel em Administração pela Universidade Estadual do Rio Grande do Sul (UERGS, 2016), Mestre (2019) e Doutor (2023) em Administração pelo Programa de Pós-Graduação em Administração da Universidade Federal de Santa Maria (UFSM). Atualmente, é Professor Doutor na Universidade Católica de Brasília e no Programa de Pós-Graduação em Governança, Tecnologia e Inovação. Desenvolve atividades de pesquisa nas áreas de economia comportamental, finanças comportamentais, mercado de capitais e análise de dados quantitativos. Sua atuação acadêmica é voltada para a compreensão do comportamento dos agentes no mercado e a utilização de métodos quantitativos para a análise de fenômenos econômicos e financeiros.

Pode ser contatado pelos e-mails: paulo.marschner@p.ucb.br

ORCID: <https://orcid.org/0000-0003-0847-2638>

LATTES: <https://lattes.cnpq.br/1245982332405570>

Prefácio

A ciência de dados consolidou-se como um dos campos mais relevantes do século XXI, impulsionada pelo crescimento exponencial da geração e do uso de dados em diferentes contextos. Nesse cenário, o e-book *Fundamentos da Ciência de Dados* apresenta-se como um guia estruturado e acessível para compreender os princípios essenciais dessa área interdisciplinar, integrando estatística, matemática, computação e domínio aplicado para a geração de conhecimento e apoio à tomada de decisões.

A obra conduz o leitor por um percurso formativo consistente, iniciando pelos conceitos fundamentais da ciência de dados e sua evolução histórica, passando por suas principais características e aplicações em setores como saúde, finanças e educação. Ao longo do texto, são apresentados os elementos centrais que sustentam a área, como o uso de grandes volumes de dados, técnicas analíticas e ferramentas tecnológicas, evidenciando seu papel estratégico na solução de problemas reais.

Destaca-se, ainda, a abordagem didática dos fundamentos estatísticos e probabilísticos, que estruturam a análise de dados. Conceitos como medidas de tendência central, dispersão, amostragem e interpretação de dados são apresentados de forma clara e aplicada, permitindo ao leitor não apenas compreender os cálculos, mas desenvolver uma leitura crítica das informações. Essa base é complementada pela introdução a temas contemporâneos, como aprendizado de máquina, inteligência artificial e visualização de dados.

Ao articular fundamentos teóricos, aplicações práticas e desafios atuais, como qualidade dos dados, ética e privacidade, o e-book vai além de um material introdutório. Ele se configura como uma base sólida para a formação em ciência de dados, promovendo uma compreensão integrada e orientada por evidências. Ao final, o leitor estará apto a interpretar dados de forma consistente e a utilizá-los como instrumento para análise e decisão em diferentes contextos.

Brasília, abril de 2026

Prof. Dr. Cláudio Chauke Nehme

Universidade Católica de Brasília

Sumário

Módulo 1

INTRODUÇÃO	12
Principais Características da Ciência de Dados.....	12
A EVOLUÇÃO DA CIÊNCIA DE DADOS	14
EXEMPLOS DE APLICAÇÃO DA CIÊNCIA DE DADOS	15
Saúde	15
Finanças.....	16
Marketing e Varejo.....	16
Indústria e Manufatura.....	17
Agricultura.....	17
Governança e Políticas Públicas	17
Educação	18
Esportes e Entretenimento	18
Sustentabilidade e Meio Ambiente.....	18
FERRAMENTAS E TECNOLOGIAS	19
Linguagens de Programação	19
Sistemas de Banco de Dados.....	19
Plataformas de Computação em Nuvem.....	19
Frameworks de Big Data.....	19
Ferramentas de Visualização de Dados.....	19
Tecnologias de Inteligência Artificial.....	20
Gestão e Automação de Fluxos de Trabalho.....	20
Integração com Internet das Coisas (IoT).....	20
Modelos Pré-treinados e APIs	20
TENDÊNCIAS	20
Aprendizagem Federada.....	21
Inteligência Artificial	21
DESAFIOS DA CIÊNCIA DE DADOS	22
Desafios Relacionados à Qualidade dos Dados	23

Desafios Relacionados à Privacidade e Segurança.....	24
Escalabilidade e Infraestrutura	25
Lacuna de Habilidades e Escassez de Talentos.....	25
Viés e Ética em Modelos	26
Manutenção e Atualização de Modelos	27
Integração de Ciência de Dados com Domínio de Negócios	27
Limitações Computacionais	28
Regulamentações e Conformidade.....	28
Sustentabilidade	29

Módulo 2

INTRODUÇÃO À PROBABILIDADE E ESTATÍSTICA BÁSICA.....	31
O que é Probabilidade.....	31
Regra de Adição em Probabilidade.....	32
Regra de Multiplicação em Probabilidade.....	33
MEDIDAS DE TENDÊNCIA CENTRAL.....	34
Média	34
Mediana.....	35
Moda	35
MEDIDAS DE DISPERSÃO.....	35
Variância	36
Desvio Padrão.....	37
POPULAÇÃO E AMOSTRA	40
População	40
Amostra	41
MÉTODOS DE AMOSTRAGEM	42
Amostragem Aleatória Simples.....	42
Amostragem Estratificada	42
Amostragem Sistemática	43
Amostragem por Conveniência.....	44
MÉDIA DA POPULAÇÃO VS. MÉDIA DA AMOSTRA	45
Média da População	45
Média da Amostra	46

VARIÁVEIS E SEUS TIPOS 47
Variáveis Qualitativas (Categóricas) 47
Variáveis Quantitativas (Numéricas). 48

DISTRIBUIÇÃO DE FREQUÊNCIA 50
Distribuição de frequência simples 50
Distribuição de frequência agrupada 51
Fórmulas para calcular frequências 51

HISTOGRAMAS 53

----- **Módulo 3**

INTRODUÇÃO À INTERPRETAÇÃO DE DADOS ESTATÍSTICOS 57

PERCENTIS E QUANTIS 57
Cálculo de Percentis 57
Uso de Tabelas para Quantis 58
Aplicações dos Percentis e Quantis 59

RESUMO DE CINCO NÚMEROS 59

INTERVALO INTERQUARTIL (IQR) 61

IDENTIFICAÇÃO DE OUTLIERS 62
Importância do IQR 62

BOXPLOTS 62
Elementos de um Boxplot 63
Visualizar a Dispersão e Identificar Outliers 63

EFEITO DE OUTLIERS E SUA REMOÇÃO 64

FUNÇÃO DENSIDADE DE PROBABILIDADE 66
Exemplos de Distribuições Comuns 66
Interpretação de Gráficos de Densidade 68

PONTUAÇÃO Z 68

Aplicações da Pontuação Z.	70
PADRONIZAÇÃO X NORMALIZAÇÃO.	70
Quando usar padronização?	71
Normalização	71
Quando usar normalização?	72
REFERÊNCIAS	73

MÓDULO 1:

Fundamentos da Ciência de Dados

Introdução

A Ciência de Dados é uma área interdisciplinar que utiliza métodos, processos, algoritmos e sistemas para extrair conhecimento e insights de dados estruturados e não estruturados. Ela combina estatística, matemática, ciência da computação e o domínio específico dos dados para solucionar problemas complexos e tomar decisões baseadas em evidências.

Figura 1 - Representação visual entre Dados Estruturados e Dados Não Estruturados.



Fonte: ChatGPT, 2026.

A Ciência de Dados pode ser definida como **o campo que transforma grandes volumes de dados em informações úteis, aproveitando técnicas como análise estatística, aprendizado de máquina e visualização de dados.** Seu escopo é amplo e abrange áreas como

- previsão de comportamentos
- identificação de padrões
- automação de processos e
- suporte à tomada de decisões estratégicas

O avanço na capacidade de coleta, armazenamento e processamento de dados, em combinação com o crescimento da computação em nuvem, possibilitou que a Ciência de Dados se expandisse rapidamente. Ela encontra aplicações em setores como saúde, finanças, marketing, indústria, segurança e muitos outros, tornando-se uma área indispensável no mundo moderno. As principais características estão listadas no Quadro 1.

Principais Características da Ciência de Dados

Característica	Descrição
Interdisciplinaridade	A Ciência de Dados combina estatística, matemática, computação e conhecimentos específicos de áreas como economia, biologia, engenharia, entre outras.
Uso de Grandes Volumes de Dados (Big Data)	Trabalha com dados em grande escala, provenientes de diversas fontes, incluindo redes sociais, dispositivos IoT, transações financeiras, entre outras.

Técnicas Analíticas	<ul style="list-style-type: none"> - Estatística e Probabilidade: Fundamentais para modelagem e análise de dados. - Aprendizado de Máquina (Machine Learning): Permite a construção de modelos preditivos e prescritivos. - Visualização de Dados: Apresentação gráfica dos resultados para facilitar a interpretação.
Automação e Escalabilidade	Ferramentas e pipelines de dados automatizados são cruciais para lidar com a escala e complexidade dos dados atuais.
Relevância Prática	A Ciência de Dados tem como objetivo principal resolver problemas reais e gerar valor por meio da informação.

Um dos principais fundamentos da Ciência de Dados está na estatística e probabilidade, que oferecem ferramentas essenciais para análise e modelagem preditiva.

- ▶ A estatística descritiva permite resumir e apresentar dados por meio de medidas como média, mediana e gráficos, enquanto a inferencial possibilita generalizações com base em amostras.
- ▶ A probabilidade serve para medir a chance de um evento ocorrer, ajudando a tomar decisões baseadas em incertezas e prever resultados possíveis em situações aleatórias.

Outro aspecto essencial é o aprendizado de máquina (Machine Learning), que permite a construção de sistemas capazes de identificar padrões e tomar decisões com base em dados. Isso inclui aprendizado supervisionado, onde modelos são treinados com dados rotulados, e não supervisionado, que identifica padrões sem rótulos predefinidos. Algoritmos como regressão linear, árvores de decisão, máquinas de vetores de suporte (SVM) e redes neurais artificiais são amplamente utilizados. O Deep Learning, uma subárea do Machine Learning, emprega redes neurais profundas para resolver problemas complexos, como visão computacional e processamento de linguagem natural, sendo viabilizado por frameworks como TensorFlow e PyTorch.



Saiba Mais

TensorFlow é uma biblioteca de software para aprendizado de máquina e inteligência artificial. Para conhecer a plataforma e saber mais sobre o assunto, acesse: <https://www.tensorflow.org/?hl=pt-br>

PyTorch é uma biblioteca de aprendizado de máquina baseada na biblioteca Torch, usada para aplicações como visão computacional e processamento de linguagem natural. Para conhecer a plataforma e saber mais sobre o assunto, acesse: <https://pytorch.org/>

A manipulação e engenharia de dados é fundamental para preparar os dados brutos antes da análise. Isso envolve etapas como limpeza, para tratar valores ausentes e inconsistentes, e transformação, como codificação de variáveis categóricas e normalização. Além disso, a integração e fusão de diferentes fontes de dados, como bancos de dados SQL e APIs, são processos importantes.

Após a preparação, a exploração e visualização de dados ajudam a identificar padrões e anomalias por meio de gráficos e dashboards, com o suporte de ferramentas como Python (Matplotlib

e Seaborn), R (ggplot2) e softwares como Tableau.

A modelagem e avaliação constituem etapas críticas, onde modelos matemáticos e estatísticos são criados para representar fenômenos ou prever eventos. A eficácia dos modelos é avaliada com métricas como erro quadrático médio (MSE) para regressão e acurácia, precisão e F1-score para classificação. Técnicas como validação cruzada e divisões treinamento-teste garantem a consistência dos modelos. Em casos de Big Data, características como volume, velocidade e variedade demandam o uso de tecnologias avançadas, como Hadoop, Spark e Kafka, para lidar com grandes volumes e alta velocidade de dados.

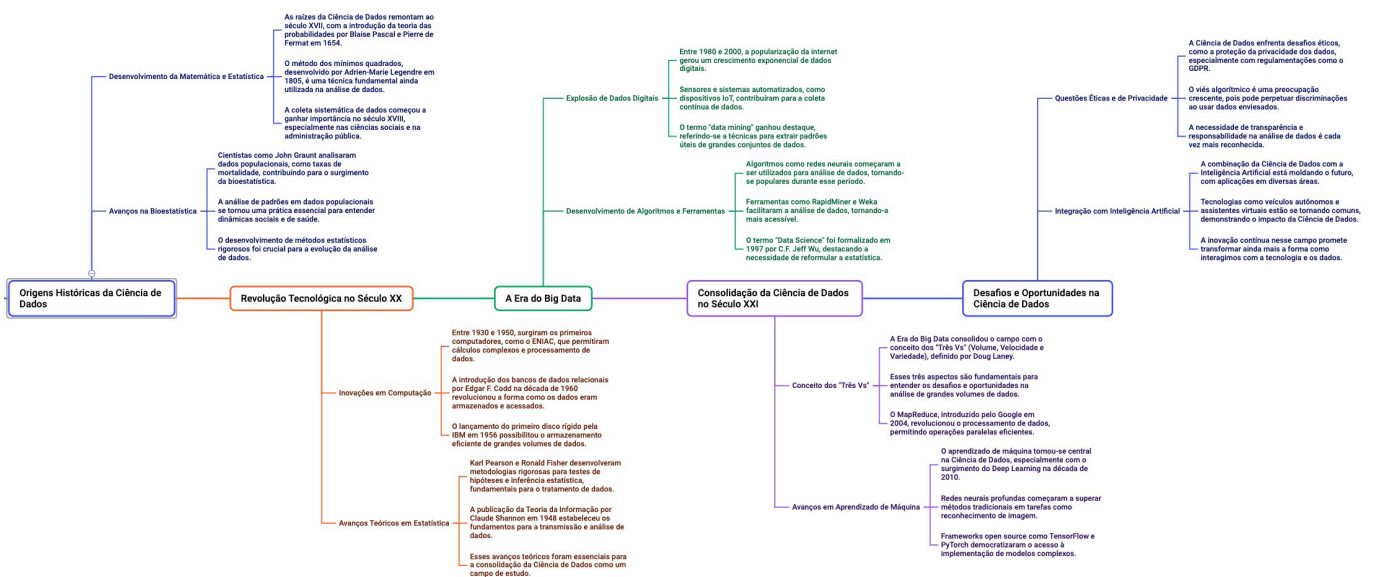
Outro pilar importante é a ética e governança de dados, que aborda questões como privacidade e viés algorítmico. Regulamentações como o GDPR e a LGPD estabelecem normas sobre o uso de dados pessoais, enquanto práticas como explicabilidade de modelos (explainable AI) são essenciais para promover confiança em áreas sensíveis. A interdisciplinaridade é outro diferencial da Ciência de Dados, combinando matemática, estatística, ciência da computação e conhecimentos específicos de diversas áreas para contextualizar insights. Essa interdisciplinaridade é fundamental para a integração com inteligência artificial, permitindo que métodos de Ciência de Dados sejam usados para treinar modelos de IA que simulam inteligência humana.

Esses aspectos-chave trabalham de forma interconectada, com fundamentos matemáticos e estatísticos formando a base teórica, enquanto a engenharia de dados prepara informações para análises avançadas. O aprendizado de máquina constrói modelos preditivos, e a visualização comunica os resultados de forma clara, apoiada por uma infraestrutura robusta para o processamento de grandes volumes de dados. Essa integração garante que a Ciência de Dados continue a transformar dados brutos em insights acionáveis, impactando setores diversos e moldando o futuro da tecnologia e da ciência.

A Evolução da Ciência de Dados

A Ciência de Dados é o resultado de séculos de avanços em matemática, estatística e computação. Apesar de ser um campo formal relativamente recente, suas raízes estão profundamente conectadas à história da ciência e tecnologia. Veja na linha do tempo a seguir a evolução da Ciência de Dados, suas raízes históricas e às suas aplicações modernas e desafios futuros.

Figura 2 - Mapa mental ciência de dados



Fonte: Elaboração própria, 2025.

O mapa mental apresentou uma visão detalhada sobre a evolução da Ciência de Dados, dividida em cinco momentos principais.

As **Origens Históricas**, destacando o desenvolvimento da matemática e estatística no século XVII, com a teoria das probabilidades e avanços na bioestatística, essenciais para análises populacionais.

A **Revolução Tecnológica no Século XX**, com o surgimento dos primeiros computadores, a introdução dos bancos de dados relacionais e os fundamentos da Teoria da Informação, que transformaram a forma de coletar e analisar dados.

A **Era do Big Data**, entre 1980 e 2000, houve uma explosão de dados digitais, impulsionada pela internet e dispositivos IoT, acompanhada pelo desenvolvimento de ferramentas como redes neurais e técnicas de "data mining".

No século XXI, a **Consolidação da Ciência de Dados** ocorre com a definição dos "Três Vs" (volume, velocidade e variedade) e a integração de tecnologias avançadas, como aprendizado de máquina e frameworks acessíveis.

Por fim, os **Desafios e Oportunidades** abordam questões éticas, como a proteção de dados e a transparência nos modelos, e ressaltam o potencial da integração com inteligência artificial para moldar o futuro do campo.

Percebemos que a evolução computacional, marcada por avanços tecnológicos significativos e a redução de custos, tornou possível a realização de análises mais complexas e acessíveis, democratizando o campo. A interdisciplinaridade desempenhou um papel central, com a convergência de estatística, matemática e computação permitindo a rápida expansão da Ciência de Dados em diferentes áreas de aplicação.

Além disso, os desafios globais contemporâneos, como as mudanças climáticas e pandemias, destacaram a relevância da Ciência de Dados como uma ferramenta indispensável para enfrentar problemas críticos.

A capacidade de extrair insights valiosos a partir de grandes volumes de dados tem sido fundamental para a tomada de decisões informadas e para o desenvolvimento de soluções inovadoras em escala global.

Exemplos de Aplicação da Ciência de Dados

A Ciência de Dados tem aplicações amplas e diversificadas, abrangendo setores como saúde, finanças, marketing, manufatura, agricultura, entre outros.

Aqui estão alguns exemplos detalhados, em como a Ciência de Dados é usada para resolver problemas reais e impulsionar a inovação.

1. Saúde

Diagnósticos e Previsões Médicas

- Exemplo: A IBM Watson utiliza aprendizado de máquina para ajudar médicos a diagnosticar doenças como câncer, analisando históricos médicos e literatura científica.
- Impacto: Redução de diagnósticos equivocados, aumento na eficiência dos tratamentos e personalização da medicina.

Predição de Epidemias

- Exemplo: Modelos preditivos baseados em dados do Google e redes sociais foram usados para monitorar a disseminação do vírus H1N1 e, mais recentemente, da COVID-19.
- Impacto: Ajudam na alocação de recursos de saúde e no planejamento de intervenções.

Análise de Imagens Médicas

- Exemplo: Ferramentas de deep learning como as desenvolvidas pela DeepMind, uma subsidiária da Google, analisam imagens de retina para prever doenças oculares.
- Impacto: Diagnósticos mais rápidos e precisos, reduzindo a dependência de especialistas humanos em áreas remotas.

2. Finanças

Prevenção de Fraudes

- Exemplo: Instituições financeiras como o PayPal utilizam algoritmos de detecção de anomalias para identificar transações fraudulentas.
- Impacto: Redução de perdas financeiras e aumento da confiança dos clientes.

Modelagem de Risco

- Exemplo: Bancos como o JPMorgan Chase empregam modelos baseados em aprendizado de máquina para prever inadimplência em empréstimos.
- Impacto: Melhor alocação de crédito e redução de perdas.

Otimização de Investimentos

- Exemplo: Robo-advisors, como o Wealthfront, usam Ciência de Dados para oferecer carteiras de investimentos personalizadas.
- Impacto: Democratização do acesso a serviços financeiros e maior eficiência na gestão de investimentos.

3. Marketing e Varejo

Personalização de Experiência do Cliente

- Exemplo: A Amazon usa modelos de recomendação baseados em comportamento de compra, histórico de navegação e preferências declaradas.
- Impacto: Aumento das vendas e da satisfação do cliente.

Análise de Sentimentos

- Exemplo: Ferramentas de análise de texto são usadas por empresas como a Coca-Cola para monitorar a percepção da marca nas redes sociais.
- Impacto: Estratégias de marketing mais eficazes e melhoria na relação com os consumidores.

Otimização de Preços

- Exemplo: Plataformas de comércio eletrônico ajustam preços dinamicamente com base em demanda, concorrência e comportamento do cliente.
- Impacto: Maximização de lucros e competitividade no mercado.

4. Indústria e Manufatura

Manutenção Preditiva

- Exemplo: A General Electric utiliza dados de sensores em equipamentos industriais para prever falhas antes que ocorram.
- Impacto: Redução de custos de manutenção e aumento da eficiência operacional.

Controle de Qualidade

- Exemplo: Empresas automotivas como a Toyota utilizam visão computacional para detectar defeitos em peças durante a produção.
- Impacto: Aumento da qualidade do produto e redução de desperdícios.

Otimização de Cadeias de Suprimento

- Exemplo: A FedEx usa modelos preditivos para planejar rotas e prever atrasos.
- Impacto: Melhor experiência para o cliente e redução de custos logísticos.

5. Agricultura

Monitoramento de Culturas

- Exemplo: Sensores IoT conectados a ferramentas de análise de dados, como os oferecidos pela John Deere, monitoram a saúde das plantações.
- Impacto: Maior produtividade agrícola e melhor uso de recursos como água e fertilizantes.

Previsão Climática

- Exemplo: Modelos baseados em aprendizado de máquina ajudam agricultores a planejar colheitas com base em previsões climáticas detalhadas.
- Impacto: Redução de perdas e melhor planejamento agrícola.

6. Governança e Políticas Públicas

Planejamento Urbano

- Exemplo: A cidade de Barcelona usa análise de dados de sensores para otimizar sistemas de transporte público e gerenciamento de resíduos.
- Impacto: Melhor qualidade de vida para os cidadãos e maior sustentabilidade urbana.

Combate ao Crime

- Exemplo: Sistemas como o PredPol preveem áreas e horários com maior probabilidade de crimes, ajudando na alocação de recursos policiais.
- Impacto: Redução de índices de criminalidade e uso mais eficiente de recursos públicos.

Análise de Políticas Públicas

- Exemplo: Governos utilizam dados para medir o impacto de programas sociais, como o Bolsa Família no Brasil.
- Impacto: Tomada de decisões baseadas em evidências e melhoria na distribuição de recursos.

7. Educação

Personalização do Ensino

- Exemplo: Plataformas como a Coursera e Khan Academy usam dados de aprendizado para personalizar o conteúdo e monitorar o progresso dos alunos.
- Impacto: Melhora na retenção de conhecimento e no engajamento dos alunos.

Previsão de Abandono Escolar

- Exemplo: Universidades utilizam modelos preditivos para identificar estudantes em risco de abandono, com base em dados de desempenho acadêmico e engajamento.
- Impacto: Redução de taxas de evasão e suporte mais direcionado aos estudantes.

8. Esportes e Entretenimento

Análise de Performance

- Exemplo: Times de basquete na NBA utilizam análise de dados para avaliar a performance de jogadores e estratégias durante os jogos.
- Impacto: Decisões mais informadas e aumento do desempenho.

Produção de Conteúdo

- Exemplo: A Netflix analisa dados de visualização para decidir quais séries e filmes produzir.
- Impacto: Melhor alinhamento entre oferta de conteúdo e preferências do público.

9. Sustentabilidade e Meio Ambiente

Monitoramento Ambiental

- Exemplo: A NASA utiliza imagens de satélite e modelos de análise de dados para monitorar o desmatamento e mudanças climáticas.
- Impacto: Políticas mais eficazes de preservação ambiental.

Gestão de Recursos Naturais

- Exemplo: Empresas de energia renovável, como a Siemens, analisam dados de geração

de energia para otimizar a distribuição em redes elétricas.

- Impacto: Maior eficiência energética e redução de desperdícios.

Ferramentas e Tecnologias

O trabalho em Ciência de Dados é viabilizado por um conjunto robusto de ferramentas e tecnologias. São elas:

1. Linguagens de Programação

As linguagens de programação são fundamentais no trabalho com dados. As principais são:

- **Python:** Reconhecida por sua versatilidade, a linguagem Python é amplamente usada devido a bibliotecas como:
 - **Pandas:** Manipulação e análise de dados.
 - **NumPy:** Processamento numérico.
 - **Scikit-learn:** Algoritmos de aprendizado de máquina.
 - **TensorFlow e PyTorch:** Modelagem de redes neurais.
- **R:** Ideal para análises estatísticas e visualizações avançadas, com pacotes como ggplot2, forecast e dplyr.

2. Sistemas de Banco de Dados

- **SQL (Structured Query Language):** Essencial para a consulta e manipulação de dados em bancos relacionais, como MySQL, PostgreSQL e Microsoft SQL Server.
- **NoSQL:** Bancos como MongoDB e Cassandra lidam melhor com dados não estruturados e de grande escala.

3. Plataformas de Computação em Nuvem

A computação em nuvem tornou a Ciência de Dados mais acessível e escalável.

- **Amazon Web Services (AWS):** Oferece ferramentas como S3 (armazenamento de dados) e SageMaker (modelagem de aprendizado de máquina).
- **Google Cloud Platform (GCP):** Inclui BigQuery para análises rápidas de grandes volumes de dados.
- **Microsoft Azure:** Suporta fluxos de trabalho de aprendizado de máquina e pipelines de dados.

4. Frameworks de Big Data

- **Apache Hadoop:** Processamento distribuído de grandes conjuntos de dados.
- **Apache Spark:** Oferece processamento em tempo real com maior velocidade e flexibilidade do que Hadoop.
- **Kafka:** Sistema de mensagens distribuídas para ingestão de dados em tempo real.

5. Ferramentas de Visualização de Dados

A visualização é fundamental para comunicar insights:

- **Tableau:** Plataforma interativa para criar dashboards.
- **Power BI:** Solução integrada da Microsoft para visualização e relatórios.
- **Matplotlib e Seaborn:** Bibliotecas Python para gráficos personalizados e análises visuais.

6. Tecnologias de Inteligência Artificial

- **Machine Learning:** Scikit-learn, TensorFlow e PyTorch são amplamente utilizados para construção de modelos preditivos.
- **Processamento de Linguagem Natural (NLP):** Ferramentas como SpaCy e NLTK ajudam a analisar textos, enquanto BERT e GPT lidam com modelos avançados de linguagem.

7. Gestão e Automação de Fluxos de Trabalho

- **Airflow:** Utilizado para orquestrar pipelines de dados complexos.
- **Luigi:** Ferramenta de automação de tarefas para pipelines menores e mais simples.

8. Integração com Internet das Coisas (IoT)

Sensores IoT produzem dados que são integrados com plataformas de análise:

- **AWS IoT Analytics:** Analisa dados provenientes de dispositivos conectados.
- **ThingSpeak:** Solução de IoT para monitoramento de dados em tempo real.

9. Modelos Pré-treinados e APIs

- **APIs da Google AI:** Reconhecimento de fala, tradução e análise de sentimentos.
- **Hugging Face:** Modelos pré-treinados de linguagem natural.

O impacto da Ciência de Dados no campo científico e social é imenso, refletindo avanços em diversas áreas do conhecimento e tecnologia. Paralelamente, o desenvolvimento constante de ferramentas e tecnologias alimenta sua evolução e acessibilidade. Essa sinergia reforça sua posição como um dos campos mais transformadores do século XXI.

Tendências

A Ciência de Dados está em constante evolução, impulsionada por avanços tecnológicos, mudanças nos volumes e tipos de dados disponíveis e novas demandas da sociedade. Entre as tendências mais relevantes está a democratização do campo, facilitada por ferramentas como AutoML (Google AutoML e H2O.ai) e plataformas no-code e low-code (Orange Data Mining, KNIME, DataRobot).

- ▶ **AutoML:** O machine learning automatizado, também conhecido como ML automatizado ou AutoML, é o processo de automatizar as tarefas demoradas e iterativas do desenvolvimento de modelo de machine learning. Com ele, cientistas de dados, analistas e desenvolvedores podem criar modelos de ML com alta escala, eficiência e produtividade, ao mesmo tempo em que dão suporte à qualidade do modelo.
- ▶ **Low-code Developer:** É uma abordagem de desenvolvimento de software que visa agilizar o processo de criação de aplicativos, permitindo que os desenvolvedores construam soluções com menos código manual e mais uso de interfaces visuais e componentes pré-construídos.



Saiba Mais

Orange Data Mining é uma ferramenta de visualização e aprendizado de máquina de código aberto. Ele possui diversos recursos para a área de Ciência de Dados. A plataforma é amigável ao usuário, versátil e tem uma variedade de aplicações.

Knime é uma plataforma de código aberto, que fornece funcionalidades como o acesso e o processamento de tipos de dados complexos, bem como a adição de algoritmos avançados de aprendizagem automática.

Essas soluções permitem que profissionais sem formação técnica avancem em análises de dados, ampliando o acesso especialmente em startups e pequenas empresas.

Aprendizagem Federada

- ▶ O aprendizado federado destaca-se como um avanço importante, possibilitando o treinamento de modelos sem centralizar dados, o que preserva a privacidade. Essa abordagem, que já é aplicada em sistemas de digitação preditiva e pesquisas médicas, atende a regulamentações como o GDPR e a LGPD, promovendo conformidade ética e segurança. Paralelamente, a expansão do Big Data, com a inclusão de dados em tempo real oriundos de dispositivos IoT e redes sociais, está transformando setores como saúde, transporte e finanças, onde a análise em tempo real é crucial. Tecnologias como Apache Kafka e Apache Flink suportam essa evolução.

Inteligência Artificial

- ▶ A integração com inteligência artificial (IA) é outro marco, com redes neurais profundas revolucionando visão computacional, processamento de linguagem natural e aplicações em saúde e marketing. Modelos como GPT e BERT estão ampliando a eficiência em resolver problemas complexos. Em paralelo, a computação quântica surge como uma tecnologia emergente com potencial de transformar a Ciência de Dados, prometendo avanços significativos em otimização e aprendizado de máquina quântico.

Monitoramento ambiental

- ▶ Questões de ética e governança de dados também ocupam o centro das atenções. Regulamentações rigorosas exigem maior transparência e responsabilidade, enquanto ferramentas como Explainable AI (XAI) ajudam a tornar algoritmos mais interpretáveis e confiáveis. Além disso, a Ciência de Dados tem sido usada para promover a sustentabilidade, com aplicações em monitoramento ambiental, agricultura inteligente e cidades conectadas. Ferramentas como o Google Earth Engine e sensores IoT estão contribuindo para otimizar recursos e mapear mudanças climáticas.



Saiba Mais

Google Earth Engine combina um catálogo de vários petabytes de imagens de satélite e conjuntos de dados geoespaciais com capacidades de análise à escala planetária. Os cientistas, investigadores e programadores utilizam o Earth Engine para detetar alterações, mapear tendências e quantificar diferenças na superfície da Terra. O Earth Engine está agora disponível para utilização comercial e permanece gratuito para utilização académica e de investigação.

Do ponto de vista de infraestrutura, arquiteturas modernas como Data Mesh, plataformas como Snowflake e Databricks, e pipelines automatizados (Apache Airflow, Prefect) estão transformando a gestão e análise de dados, promovendo eficiência e escalabilidade. Essas inovações têm impacto direto em setores como saúde, finanças e marketing, que se beneficiam de modelos preditivos, personalização de serviços e análise comportamental.



Saiba Mais

Databricks é uma empresa de dados e IA. O seu site discute a plataforma, os produtos e os recursos da empresa. A plataforma Databricks unifica dados, análises e IA. A empresa oferece soluções para vários setores, incluindo saúde, varejo e manufatura. A Databricks também oferece treinamento, certificações e eventos.

Olhando para o futuro, a integração de dados multimodais será uma área de destaque, permitindo análises simultâneas de textos, imagens e áudios. Essa tecnologia terá impacto em setores como saúde, com análises clínicas avançadas, e educação, com ferramentas interativas de aprendizado. Modelos de Explainable AI (XAI) continuarão a evoluir, atendendo à demanda por transparência e regulamentações mais rigorosas.

A computação quântica e o aprendizado federado prometem transformações profundas, enquanto a automação da Ciência de Dados continuará a democratizar o acesso, tornando-a uma ferramenta essencial para pequenas e médias empresas. O foco crescente em sustentabilidade levará à criação de modelos energeticamente eficientes e ao uso de infraestrutura alimentada por energia renovável.

A Ciência de Dados continuará a ampliar sua relevância técnica, social e econômica, moldando um futuro em que privacidade, transparência e sustentabilidade serão pilares fundamentais para a inovação e o progresso global.

Desafios da Ciência de Dados

A Ciência de Dados enfrenta desafios técnicos, éticos e organizacionais que podem limitar seu impacto e crescimento. No entanto, soluções inovadoras, como novas arquiteturas, maior transparência e regulamentações robustas, estão ajudando a superar essas barreiras. À medida que a área evolui, a colaboração interdisciplinar e o compromisso com a ética serão fundamentais para garantir que a Ciência de Dados continue a transformar positivamente a sociedade.

A seguir vamos conhecer os principais desafios da área, detalhando suas causas, impactos e possíveis soluções.

1. Desafios Relacionados à Qualidade dos Dados

Dados Não Estruturados e Heterogeneidade

Grande parte dos dados disponíveis atualmente é não estruturada, como imagens, vídeos, áudios e textos. Trabalhar com esses dados apresenta desafios como:

Formato e Organização:

Dados não estruturados não seguem um esquema fixo, tornando sua análise mais complexa.

Heterogeneidade:

Dados de diferentes fontes (sensores IoT, redes sociais, sistemas financeiros) apresentam formatos e padrões variados, exigindo esforços de integração e normalização.

Dados Ausentes e Ruidosos

Dados incompletos ou com ruído comprometem a precisão dos modelos:

Impacto:

Modelos baseados em dados inconsistentes podem produzir resultados imprecisos, levando a decisões equivocadas com consequências significativas.

Exemplo:

Em sistemas de saúde, dados ausentes sobre pacientes podem prejudicar a personalização de tratamentos.

Soluções Propostas

Técnicas de Imputação:

Métodos estatísticos para preencher valores ausentes.

Pipelines Automatizados:

Ferramentas como Apache Airflow e Prefect auxiliam na limpeza e organização dos dados.

Deteção e Remoção de dados Ruidosos:

Técnicas de deteção de anomalias, como o uso de algoritmos de isolamento, podem ser aplicadas para identificar e eliminar dados ruidosos, melhorando a qualidade dos dados.

2. Desafios Relacionados à Privacidade e Segurança

Preocupações com Privacidade

A coleta e análise de grandes volumes de dados levantam questões éticas:

Dados Pessoais Sensíveis:

Informações como localização, saúde e histórico financeiro podem ser exploradas de forma inadequada.

Regulamentações:

Leis como a GDPR (União Europeia) e a LGPD (Brasil) impõem restrições sobre o uso e armazenamento de dados pessoais.

Segurança de Dados

Ataques cibernéticos e violações de segurança comprometem a integridade dos dados:

Exemplo:

Vazamentos em plataformas como Facebook e Equifax expuseram milhões de usuários a riscos de roubo de identidade.

Soluções Propostas

Criptografia e Anonimização:

Métodos para proteger dados sensíveis durante a coleta e armazenamento.

Aprendizado Federado:

Permite treinar modelos sem a necessidade de compartilhar os dados brutos.

Autenticação Multifatorial:

O uso de autenticação multifatorial oferece uma camada extra de segurança para o acesso a sistemas e dados sensíveis, protegendo contra acessos não autorizados em caso de roubo de credenciais.

3. Escalabilidade e Infraestrutura

Volume e Velocidade de Dados

A era do Big Data trouxe desafios relacionados ao processamento de grandes volumes de dados em tempo real:

Crescente Volume de Dados:

Sensores IoT, redes sociais e transações financeiras geram dados em escalas sem precedentes.

Impacto na Infraestrutura:

Sistemas precisam lidar com armazenamento, processamento e análise em tempo real.

Soluções Propostas

Arquiteturas de Big Data:

Tecnologias como Apache Hadoop e Spark permitem o processamento distribuído.

Computação em Nuvem:

Plataformas como AWS e Google Cloud oferecem escalabilidade sob demanda.

4. Lacuna de Habilidades e Escassez de Talentos

Escassez de Cientistas de Dados

A demanda por profissionais qualificados em Ciência de Dados supera a oferta, resultando em uma lacuna de habilidades no mercado:

Impacto:

Muitas organizações não conseguem implementar soluções de dados eficazes por falta de profissionais capacitados.

Dificuldade de Formação

A área exige conhecimentos multidisciplinares, incluindo estatística, programação, aprendizado de máquina e domínio do negócio:

Barreira de Entrada:

A necessidade de dominar várias ferramentas e linguagens cria dificuldades para novos profissionais.

Soluções Propostas

Educação e Treinamento:

Cursos online e programas de certificação, como os oferecidos por plataformas como Coursera e Udemy, ajudam a formar novos profissionais.

Ferramentas de No-Code:

Simplificam processos, permitindo que profissionais de outras áreas utilizem ciência de dados.

5. Viés e Ética em Modelos

Viés nos Dados e Modelos

Dados enviesados podem levar a decisões injustas ou discriminatórias:

Exemplo:

Algoritmos de contratação treinados em dados históricos podem perpetuar desigualdades de gênero ou raça.

Falta de Transparência

Modelos de aprendizado profundo são frequentemente criticados por serem "caixas pretas", dificultando a explicação de suas decisões:

Impacto:

Em setores como saúde e finanças, a falta de transparência pode comprometer a confiança nos sistemas.

Soluções Propostas

Mitigação de Viés:

Ferramentas como IBM AI Fairness 360 ajudam a identificar e corrigir vieses.

Explainable AI (XAI):

Técnicas como SHAP e LIME tornam os modelos mais interpretáveis.

6. Manutenção e Atualização de Modelos

Obsolescência de Modelos

Modelos de Ciência de Dados precisam ser atualizados regularmente para se manterem precisos:

Drift de Dados:

Mudanças no comportamento dos dados ao longo do tempo podem reduzir a eficácia dos modelos.

Soluções Propostas

Monitoramento Contínuo:

Ferramentas para avaliar o desempenho dos modelos em tempo real.

Aprendizado Contínuo:

Métodos que permitem que os modelos sejam ajustados constantemente.

7. Integração de Ciência de Dados com Domínio de Negócios

Alinhamento com Objetivos

Muitas vezes, as soluções de dados não estão alinhadas com os objetivos estratégicos das organizações:

Impacto:

Projetos de dados podem consumir recursos significativos sem gerar valor tangível.

Comunicação de Resultados

Os cientistas de dados enfrentam desafios para comunicar resultados complexos a públicos não técnicos:

Exemplo:

Executivos podem não compreender totalmente as limitações e incertezas de um modelo.

Soluções Propostas

Colaboração Interdisciplinar:

Envolver especialistas de domínio no desenvolvimento de projetos.

Visualização de Dados:

Ferramentas como Tableau e Power BI facilitam a comunicação de insights.

8. Limitações Computacionais

Modelos de Alta Complexidade

Algoritmos de aprendizado profundo, como redes neurais profundas, demandam grande poder computacional:

Impacto:

O custo e o tempo de processamento podem ser proibitivos para pequenas empresas ou organizações.

Soluções Propostas

Uso de GPUs e TPUs:

Hardware especializado acelera o treinamento de modelos.

Computação Quântica:

Embora ainda em desenvolvimento, promete revolucionar a análise de dados complexos.

9. Regulamentações e Conformidade

Barreiras Legais

A conformidade com regulamentações como LGPD e GDPR pode limitar o acesso e uso de dados:

Exemplo:

Empresas que operam em múltiplas jurisdições precisam atender a diferentes regulamentações, aumentando a complexidade operacional.

Soluções Propostas

Compliance Integrado:

Ferramentas que automatizam verificações de conformidade.

Governança de Dados:

Estratégias para gerenciar dados de forma ética e legal.

10. Sustentabilidade

Consumo de Energia

Modelos de grande escala, como os usados em aprendizado profundo, consomem quantidades significativas de energia:

Impacto Ambiental:

O treinamento de um único modelo pode emitir tanto carbono quanto um carro em um ano.

Soluções Propostas

Otimização de Algoritmos:

Métodos mais eficientes para reduzir o consumo de energia.

Uso de Fontes Renováveis:

Data centers alimentados por energia solar ou eólica.

MÓDULO 2:

*Explorando a Estatística:
Probabilidade, Dados e Análises*

Introdução à Probabilidade e Estatística Básica

Neste módulo, apresentaremos o campo da probabilidade, uma área da matemática que tem o poder de medir e quantificar a chance de um evento acontecer.

Figura 3 - Jogo de dados



Fonte: SHVETS production, Pexels 2025.

Esse conceito, que permeia nosso cotidiano, é **fundamental para entender como decisões são tomadas em cenários de incerteza**. Desde escolhas aparentemente simples, como decidir se devemos ou não levar um guarda-chuva ao sair de casa, até questões mais complexas, como a previsão de comportamentos de mercado, a probabilidade está sempre presente.

Além de ser uma ferramenta matemática, a probabilidade também atua como um elo entre várias disciplinas.

- Na economia, por exemplo, ela é usada para modelar riscos e oportunidades;
- Na biologia, para prever a distribuição de características genéticas nas populações, com base em modelos probabilísticos;
- Na ciência da computação, para otimizar algoritmos que aprendem com dados.

Ao mergulharmos no estudo da probabilidade, veremos como esse conceito se entrelaça com a análise de grandes volumes de dados, oferecendo uma base sólida para a interpretação de tendências e a previsão de resultados futuros.

Em um mundo cada vez mais dependente da ciência de dados, o domínio da probabilidade nos habilita a fazer inferências mais precisas, seja em estudos científicos, na análise de comportamentos humanos ou até no desenvolvimento de soluções tecnológicas inovadoras.

O que é Probabilidade

A essência da probabilidade reside em quantificar o grau de incerteza de eventos e usá-lo para prever a probabilidade de diferentes desfechos. Em termos matemáticos, a probabilidade de um evento A acontecer é dada pela razão entre o número de resultados favoráveis a esse evento e o número total de resultados possíveis.

A fórmula é expressa como:

$$P(A) = \frac{\text{número de resultados favoráveis a } A}{\text{número total de resultados possíveis}}$$

Figura 4 - Dados empilhados

Para ilustrar esse conceito, considere o exemplo clássico de um dado de seis faces. Se quisermos saber qual é a probabilidade de obter um número par ao lançá-lo, precisamos primeiro identificar quantos números pares existem entre as possíveis opções (1, 2, 3, 4, 5, 6). Neste caso, os números pares são 2, 4 e 6, ou seja, três resultados favoráveis. Como o total de resultados possíveis é 6, a probabilidade de sair um número par é:

$$P(\text{par}) = \frac{3}{6} = 0,5$$



Fonte: wirestock, Freepik, 2026.

Isso significa que, a cada vez que jogamos o dado, temos 50% de chance de obter um número par. Esse princípio é a base de muitas aplicações da probabilidade, permitindo a análise de situações aleatórias e ajudando a interpretar padrões de forma precisa.

Seja na previsão do tempo, na análise de dados genéticos ou na determinação de padrões de comportamento em grandes volumes de dados, a probabilidade oferece uma ferramenta poderosa para compreender o mundo ao nosso redor.

Regra de Adição em Probabilidade

A regra de adição é utilizada em probabilidade quando queremos calcular a chance de que um de vários eventos ocorra. Esse princípio é particularmente útil quando estamos lidando com situações onde existem múltiplos resultados possíveis e queremos entender a probabilidade de ocorrência de pelo menos um desses resultados. A regra de adição é frequentemente aplicada a eventos que são mutuamente exclusivos, ou seja, eventos que não podem ocorrer simultaneamente. Para eventos mutuamente exclusivos, a probabilidade de que o evento A ou o evento B ocorra é simplesmente a soma das probabilidades de cada evento individual.

A fórmula é expressa como:

$$P(A \text{ ou } B) = P(A) + P(B)$$

Por exemplo, vamos considerar um baralho padrão de 52 cartas.

Suponha que desejamos calcular a probabilidade de sacar um **Rei** ou uma **Rainha** em uma única tentativa. Como não é possível tirar ao mesmo tempo um Rei e uma Rainha no mesmo saque, esses eventos são mutuamente exclusivos. Sabemos que existem 4 Reis e 4 Rainhas em um baralho. Portanto, a probabilidade de sacar um Rei é:

$$P(\text{Rei}) = \frac{4}{52}$$

Da mesma forma, a probabilidade de sacar uma Rainha é:

$$P(\text{Rainha}) = \frac{4}{52}$$

Para calcular a probabilidade de sacar **um Rei ou uma Rainha**, aplicamos a regra de adição:

$$P(\text{Rei ou Rainha}) = \frac{4}{8} + \frac{4}{8} = \frac{8}{52}$$

Essa fração pode ser simplificada para $\frac{2}{13}$, ou aproximadamente 15,38%. Isso significa que em um saque de uma única carta, a chance de obter um Rei ou uma Rainha é cerca de 15,38%. A regra de adição é amplamente utilizada em ciência de dados por combinar probabilidades de diferentes cenários para prever padrões e fazer inferências, seja em relação a comportamento humano, resultados econômicos ou padrões naturais. Assim, ao dominar essa regra, você amplia sua capacidade de trabalhar com incertezas de maneira sistemática e fundamentada.

Regra de Multiplicação em Probabilidade

A regra de multiplicação é utilizada quando queremos calcular a chance de que dois ou mais eventos ocorram simultaneamente. Ela é aplicada quando estamos lidando com eventos independentes, ou seja, eventos em que o resultado de um não afeta o resultado do outro. A regra de multiplicação é particularmente útil para entender a probabilidade conjunta de eventos que ocorrem em sequência ou ao mesmo tempo, sendo amplamente usada em áreas como estatística, ciência de dados, genética e diversas outras disciplinas que dependem de cálculos probabilísticos. Para eventos independentes, a probabilidade de que o evento A e o evento B ocorram simultaneamente é o produto das probabilidades de cada evento.

A fórmula é expressa como:

$$P(A \text{ e } B) = P(A) \times P(B)$$

Para ilustrar esse conceito, imagine que estamos jogando dois dados de seis faces.

Qual seria a probabilidade de que ambos mostrem o número 6 ao mesmo tempo? Como o resultado de um dado não influencia o outro, esses eventos são independentes. Sabemos que a probabilidade de que um dado mostre o número 6 é:

$$P(6) = \frac{1}{6}$$

Como os eventos são independentes, aplicamos a regra de multiplicação para calcular a probabilidade de que ambos os dados mostrem 6:

$$P(6 \text{ e } 6) = \frac{1}{6} \times \frac{1}{6} = \frac{1}{36}$$

Portanto, a probabilidade de que ambos os dados mostrem o número 6 é $\frac{1}{36}$, ou aproximadamente 2,78%. Esse tipo de cálculo é extremamente útil em diversas situações da vida real e em diferentes áreas de estudo.

Na ciência de dados, a regra de multiplicação é aplicada para modelar a probabilidade de eventos simultâneos em grandes conjuntos de dados, como a probabilidade de um sistema falhar em diferentes componentes ao mesmo tempo ou a chance de um determinado padrão se repetir em diferentes condições.

Medidas de Tendência Central

As **Medidas de Tendência Central** são ferramentas estatísticas que representam o ponto central ou típico de um conjunto de dados. Elas ajudam a resumir um grande volume de informações em um único valor, que reflete o comportamento mais comum ou esperado dos dados.

► As três principais medidas de tendência central são: **média, mediana e moda.**

Essas ferramentas são essenciais para resumir e descrever um conjunto de dados de maneira clara e compreensível. Ao identificar a "posição central" de um grupo de valores, essas medidas fornecem uma visão do comportamento típico ou padrão dos dados, permitindo que façamos inferências e tomemos decisões informadas.

Medidas de tendência central são amplamente usadas em diversas disciplinas, como ciência de dados, estatística, economia, saúde e ciências sociais, além de terem aplicações práticas em situações cotidianas, como avaliar notas de alunos, calcular médias de renda e até prever tendências de consumo.

Vamos entender melhor cada uma.

1. Média

A média é talvez a medida de tendência central mais conhecida. Ela é calculada somando-se todos os valores de um conjunto de dados e dividindo o resultado pelo número total de valores. A média é uma excelente forma de resumir os dados, fornecendo uma estimativa do "valor típico" de um conjunto.

No entanto, é importante lembrar que ela pode ser influenciada por valores extremos (*outliers*).

A fórmula para calcular a média é:

$$\bar{x} = \frac{\sum x}{n}$$

Onde:

$\sum x$ é a soma de todos os valores do conjunto de dados.

n é o número total de valores.

Exemplo: Imagine que temos as idades de cinco pessoas: 18, 20, 22, 19 e 21 anos. A média dessas idades é:

$$\bar{x} = \frac{18 + 20 + 22 + 19 + 21}{5} = 20$$

Isso significa que a idade média das cinco pessoas é 20 anos.

2. Mediana

A mediana é o valor que se encontra no meio de um conjunto de dados quando eles são organizados em ordem crescente. Diferente da média, a mediana não é afetada por valores extremos, o que a torna uma medida mais robusta para representar a tendência central em conjuntos de dados que possuem *outliers* ou distribuições assimétricas.

Para encontrar a mediana:

- Se o número de valores for ímpar, a mediana é o valor central.
- Se o número de valores for par, a mediana é a média dos dois valores centrais.

Exemplo: Se temos os seguintes valores: 10, 15 e 20, ao organizá-los em ordem crescente, vemos que o valor central é 15. Portanto, a mediana é 15.

3. Moda

A moda é o valor que aparece com maior frequência em um conjunto de dados. É especialmente útil quando estamos interessados em identificar o valor mais comum ou frequente em um conjunto de observações. A moda pode ser única, ou um conjunto de dados pode ter mais de uma moda (multimodal), ou ainda não ter nenhuma moda, caso nenhum valor se repita.

Exemplo: Considere os valores 2, 3, 4, 3 e 5. O número 3 aparece duas vezes, enquanto os outros valores aparecem apenas uma vez. Portanto, a moda desse conjunto é 3.

Essas três medidas — média, mediana e moda — formam a base para a compreensão de conjuntos de dados e são ferramentas essenciais na ciência de dados.

Agora que você já entendeu que **média, mediana e moda** indicam o ponto central ou valor representativo de um conjunto de dados, vamos ver o quão espalhados ou concentrados os dados estão ao redor dessas medidas centrais.

Para precisarmos compreender algumas medidas de dispersão.

Medidas de Dispersão

Enquanto as medidas de tendência central, como a média, a mediana e a moda, nos fornecem uma visão geral da posição central dos dados, as **medidas de dispersão** nos ajudam a entender o quão espalhados ou concentrados esses dados estão em torno dessa posição central. Saber como os dados se distribuem nos oferece uma compreensão mais detalhada da consistência dos dados e dos possíveis padrões de variação. Isso é essencial para prever com maior precisão resultados futuros e identificar comportamentos atípicos ou extremos nos dados.

Figura 5 - Charge Pere Roca



Fonte: Pere Roca, Psicometria, 2023.

Medidas de dispersão são amplamente utilizadas em muitas áreas do conhecimento, como estatística, economia, saúde, engenharia e até em estudos de comportamento humano. Elas fornecem uma base sólida para a análise de dados, especialmente quando precisamos avaliar o grau de incerteza ou variabilidade em um conjunto de informações.

Variância

A variância é uma medida que indica o quão dispersos ou afastados os valores de um conjunto de dados estão em relação à sua média. Ela calcula a média dos quadrados das diferenças entre cada valor e a média do conjunto, oferecendo uma visão numérica da extensão da variabilidade dos dados. Quanto maior a variância, mais espalhados estão os valores em relação à média. Valores baixos de variância indicam que os dados estão mais concentrados em torno da média.

A fórmula para calcular a variância é:

$$\sigma^2 = \frac{\sum(x - \bar{x})^2}{n}$$

Onde:

σ^2 é a variância

x representa cada valor no conjunto de dados.

\bar{x} é a média dos valores.

n é o número total de valores.

Exemplo: Vamos calcular a variância dos seguintes valores: 2, 4, 4, 6 e 8.

Primeiro, calculamos a média:

$$\bar{x} = \frac{2 + 4 + 4 + 6 + 8}{5} = \frac{24}{5} = 4,8$$

Em seguida, calculamos a diferença de cada valor em relação à média, elevamos ao quadrado e somamos:

$$(2 - 4,8)^2 = 7,87$$

$$(4 - 4,8)^2 = 0,64$$

$$(4 - 4,8)^2 = 0,64$$

$$(6 - 4,8)^2 = 1,44$$

$$(8 - 4,8)^2 = 10,24$$

Soma das diferenças ao quadrado: $7,84 + 0,64 + 0,64 + 1,44 + 10,24 = 20,8$.

Finalmente, dividimos pela quantidade de valores para encontrar a variância:

$$\sigma^2 = \frac{20,8}{5} = 4,16$$

Portanto, a variância desse conjunto de dados é 4,16.

A **variância** é amplamente utilizada em diversas situações reais para medir a dispersão ou variabilidade dos dados em torno de uma média. Veja alguns exemplos:

Finanças

Risco de Investimentos:

Em análises de mercado, a variância é usada para medir a volatilidade de ações ou portfólios. Um investimento com alta variância indica maior risco, pois os retornos esperados podem variar significativamente da média.

Saúde

Estudos Epidemiológicos:

A variância ajuda a entender a dispersão de variáveis como pressão arterial, níveis de colesterol ou peso em uma população, permitindo identificar padrões ou grupos de risco.

Climatologia

Estudos de Variação Climática:

Cientistas utilizam a variância para analisar as flutuações de temperatura, precipitação ou outras condições meteorológicas ao longo de períodos.

Percebeu como a variância é uma ferramenta essencial para medir a consistência e entender padrões em diferentes contextos?

E por falar em padrões, vamos entender como eles funcionam.

Desvio Padrão

Imagine que você começou a monitorar seus gastos diários para entender melhor suas finanças. E pra começar, você quer saber **se está gastando de forma equilibrada no dia a dia**.

Em uma semana, você registrou valores como 50, 52, 70, 49, 48, 53 e 75 reais. Olhando rapidamente, você pode pensar: *"Será que estou sendo consistente ou há algo fora do padrão?"*

As possibilidades são várias: talvez alguns dias exijam gastos inesperados, como compras maiores ou contas imprevistas, enquanto outros sejam mais controlados.

Mas como ter certeza de que esses desvios são significativos e não apenas uma percepção?

O **desvio padrão** entra aqui como a solução para medir o quanto seus gastos variam em relação à média semanal. Se a variação for alta, isso indica que você não está tão regular quanto imagina e pode ser hora de ajustar seus hábitos para alcançar maior equilíbrio. Se for baixa, é um sinal positivo de que você tem controle financeiro. Essa análise simples pode ser a chave para organizar suas finanças de forma mais eficiente!

O desvio padrão é uma medida de dispersão que indica, em média, o quanto os valores de um conjunto de dados se desviam da média. Ele é simplesmente a raiz quadrada da variância e oferece uma interpretação mais intuitiva da dispersão dos dados, já que retorna os desvios na mesma unidade dos dados originais. O desvio padrão é amplamente utilizado em análises estatísticas para medir a variabilidade e a consistência de um conjunto de dados.

A fórmula para calcular o desvio padrão é:

$$\sigma = \sqrt{\frac{\sum(x - \bar{x})^2}{n}}$$

Exemplo: Usando os valores da variância do exemplo anterior (4,16), o desvio padrão é:

$$\sigma = \sqrt{4,16} = 2,04$$

Isso significa que, em média, os valores se desviam da média em aproximadamente 2,04 unidades.

Entender a variância e o desvio padrão permite uma visão mais detalhada sobre como os dados se distribuem ao redor da média. Esses conceitos são aplicáveis em praticamente todas as áreas do conhecimento e são fundamentais para a ciência de dados, onde a precisão nas análises depende diretamente de como tratamos a dispersão dos dados.

Agora que você já viu o que é desvio padrão, vamos te mostrar na prática como fazer os cálculos, por meio da situação apresentada no início desse tópico.

Vamos calcular o **desvio padrão** para os gastos diários: **50, 52, 70, 49, 48, 53, 75 reais**.

Passo 1: Calcular a média (\bar{x})

Somando todos os valores e dividimos pelo número de dias:

$$\text{Média} = \frac{50 + 52 + 70 + 49 + 48 + 53 + 75}{7} = \frac{397}{7} = 56,71 \text{ reais}$$

Passo 2: Calcular as diferenças em relação à média ($x - \bar{x}$)

- Para 50: $50 - 56,71 = -6,71$
- Para 52: $52 - 56,71 = -4,71$
- Para 70: $70 - 56,71 = 13,29$
- Para 49: $49 - 56,71 = -7,71$
- Para 48: $48 - 56,71 = -8,71$
- Para 53: $53 - 56,71 = -3,71$
- Para 75: $75 - 56,71 = 18,29$

Passo 3: Elevar as diferenças ao quadrado ($(x - \bar{x})^2$)

Para cada valor subtraímos a média:

- Para 50: $(-6,71)^2 = 45,03$
- Para 52: $(-4,71)^2 = 22,18$
- Para 70: $(13,29)^2 = 176,62$
- Para 49: $(-7,71)^2 = 59,45$
- Para 48: $(-8,71)^2 = 75,89$
- Para 53: $(-3,71)^2 = 13,76$
- Para 75: $(18,29)^2 = 334,54$

Passo 4: Calcular a variância (σ^2)

Somamos os valores obtidos e dividimos pelo número de dias

$$\sigma^2 = \frac{45,03 + 22,18 + 176,62 + 59,45 + 75,89 + 13,76 + 334,54}{7} = \frac{727,47}{7} \approx 103,92$$

Passo 5: Calcular o desvio padrão (σ)

O desvio padrão é a raiz quadrada da variância:

$$\sigma = \sqrt{103,92} \approx 10,19 \text{ reais}$$

Resultado:

O desvio padrão é **10,19 reais**, indicando que os gastos diários variam, em média, cerca de **10,19 reais** em relação à média de **56,71 reais**. Esse valor relativamente alto mostra que há inconsistência nos gastos ao longo da semana, com dias muito acima ou abaixo do padrão. Isso sugere a necessidade de avaliar os dias com gastos fora do comum (como 70 e 75 reais) para identificar possíveis ajustes no controle financeiro.

População e Amostra

Se você quisesse saber qual é o melhor horário para todos os seus amigos se reunirem, seria possível perguntar a cada um deles? Ou bastaria perguntar a um pequeno grupo? Como garantir que esse grupo realmente representa a escolha de todos?

Essa pergunta nos leva a refletir sobre o conceito de **população e amostra**. A população representa o grupo completo que queremos entender, como todos os seus amigos, enquanto a amostra é uma parte menor desse grupo que usamos para tirar conclusões. Para que as respostas de um pequeno grupo sejam representativas da população, é essencial escolher a amostra de forma cuidadosa, garantindo que ela reflita as características do grupo maior. Assim, podemos economizar tempo e esforço ao mesmo tempo em que obtemos informações precisas para tomar decisões.

Enquanto a **população** representa o conjunto total de elementos em estudo, a **amostra** é um subconjunto representativo usado para fazer inferências.

Este conhecimento é a base para a coleta de dados e análise estatística, possibilitando que tiremos conclusões sobre grandes populações sem precisar estudá-las integralmente.

Então vamos nos aprofundar melhor nesses conceitos?

População

A **população** é o conjunto completo de todos os elementos ou indivíduos que estão sendo estudados ou sobre os quais queremos tirar conclusões. Cada indivíduo ou elemento da população tem uma característica em comum, que é o foco do estudo. A população pode ser grande, como todos os habitantes de um país, ou menor, como todos os alunos de uma sala de aula.

Exemplos de População



Saúde

Se um pesquisador quer estudar a taxa de infecção por uma doença em um hospital, a **população** seria formada por todos os pacientes que passaram pelo hospital em um determinado período.



Economia

Ao estudar o consumo de energia elétrica em uma cidade, a **população** seria composta por todas as residências da cidade que utilizam eletricidade.



Agricultura

Se estamos interessados em analisar a produtividade de uma safra de soja em uma fazenda, a **população** seria composta por todas as plantas de soja daquela fazenda.

No dia a dia, entender o conceito de população é importante porque, em muitos estudos, não é viável coletar dados de todos os indivíduos de uma população devido às limitações de tempo, custo ou acesso. Por isso, muitas vezes recorreremos à amostragem.

Amostra

A **amostra** é um subconjunto da população que é selecionado para ser estudado. A ideia é que a amostra seja representativa da população, ou seja, que os dados coletados a partir da amostra possam ser usados para fazer inferências sobre a população como um todo. Ao escolhermos uma amostra adequadamente, podemos economizar recursos e ainda assim obter resultados significativos.

Exemplos de Amostra



Saúde

Em vez de testar todos os pacientes de um hospital para medir a taxa de infecção por uma doença, os pesquisadores podem coletar dados de uma **amostra** de 100 pacientes escolhidos aleatoriamente.



Economia

Se quisermos saber o gasto médio de eletricidade em uma cidade, podemos selecionar uma **amostra** de 500 casas em vez de estudar todas as residências.



Agricultura

Para avaliar a produtividade de uma plantação de soja, podemos selecionar uma **amostra** de 100 plantas de soja de diferentes áreas da fazenda, em vez de medir todas as plantas.

A amostra precisa ser representativa da população para garantir que as conclusões obtidas sejam válidas. Por exemplo, se uma pesquisa de consumo de eletricidade em uma cidade incluir apenas casas de bairros ricos, a média obtida pode não representar corretamente o consumo de eletricidade das áreas mais pobres, e a análise estará enviesada.

Para não confundir os termos, veja na tabela a seguir a diferença entre cada um.

Característica	População	Amostra
Definição	Conjunto completo de todos os indivíduos ou elementos em estudo.	Subconjunto da população selecionado para estudo.
Tamanho	Grande, pode incluir todos os elementos de interesse (ex.: todos os estudantes de uma escola).	Menor, inclui apenas parte dos elementos da população (ex.: 50 estudantes de uma escola).
Uso	Usada quando é possível ou necessário estudar todos os elementos.	Usada quando não é viável ou prático estudar toda a população.
Objetivo	Obter informações precisas e completas sobre o grupo total.	Fazer inferências ou estimativas sobre a população a partir do estudo de uma parte.

Custo e Tempo	Maior custo e mais tempo necessário para coletar dados de toda a população.	Menor custo e tempo para coletar dados de uma amostra representativa.
Importância da Representatividade	Não se aplica, pois todos os elementos estão incluídos.	A amostra deve ser representativa para garantir que os resultados possam ser generalizados.

Métodos de Amostragem

Os **métodos de amostragem** são formas de selecionar uma parte representativa da população para análise. Escolher corretamente o método de amostragem é essencial para garantir que a amostra represente a população de forma precisa, permitindo a obtenção de dados confiáveis e conclusões válidas.

Amostragem Aleatória Simples

A **amostragem aleatória simples** é um método no qual cada membro da população tem a mesma chance de ser escolhido. Esse método é o mais básico e amplamente utilizado, pois garante que a seleção seja justa e sem viés.

Exemplos do Cotidiano



Saúde

Um hospital deseja estudar a satisfação dos pacientes. Com esse método, cada paciente que passou pelo hospital tem a mesma probabilidade de ser selecionado para participar da pesquisa.



Economia

Para entender os padrões de consumo em uma cidade, uma empresa pode selecionar aleatoriamente um número de domicílios para uma pesquisa de consumo mensal.



Agricultura

Se quisermos estudar o rendimento de um campo de milho, podemos escolher aleatoriamente plantas em diferentes partes do campo para medir sua produtividade.

Essa técnica é simples e evita viés de seleção, garantindo que todos os membros da população tenham a mesma chance de serem escolhidos.

Amostragem Estratificada

Na **amostragem estratificada**, a população é dividida em subgrupos, chamados **estratos**, com base em uma característica específica, como idade, gênero, região, ou categoria profissional. Em seguida, seleciona-se uma amostra aleatória de cada estrato. Esse método é útil quando a população é heterogênea e há subgrupos com características distintas.

Exemplos do Cotidiano



Saúde

Para estudar o impacto de um novo medicamento, os pacientes podem ser divididos em estratos com base na idade. Em seguida, uma amostra proporcional de cada faixa etária é selecionada para garantir que o estudo reflita a diversidade de idade.



Economia

Em uma pesquisa de renda familiar, as famílias podem ser divididas em estratos com base na renda anual, e uma amostra proporcional de cada faixa de renda é coletada para análise.



Agricultura

Um pesquisador quer entender a produção de diferentes tipos de frutas em uma fazenda. A fazenda é dividida em estratos de acordo com os tipos de frutas plantadas (maçã, laranja, banana), e amostras são coletadas de cada estrato.

Esse método garante que todas as subcategorias da população sejam representadas, o que torna os resultados mais precisos e aplicáveis a diferentes grupos dentro da população.

Amostragem Sistemática

Na **amostragem sistemática**, os indivíduos são selecionados de maneira regular, ou seja, o pesquisador escolhe o primeiro indivíduo aleatoriamente e depois seleciona outros de acordo com um intervalo fixo. Esse método é eficiente e mais fácil de aplicar do que a amostragem aleatória simples.

Exemplos do Cotidiano



Saúde

Em um hospital, se houver 1.000 pacientes registrados, pode-se selecionar a cada 10º paciente para participar de uma pesquisa sobre a qualidade do atendimento.



Economia

Uma empresa quer analisar as compras feitas em uma loja. Eles selecionam o primeiro cliente do dia e, em seguida, entrevistam cada 5º cliente que passa pelo caixa.



Agricultura

Ao medir a qualidade do solo em uma plantação, o pesquisador escolhe um ponto inicial e, a cada 20 metros, coleta uma amostra de solo para análise.

É mais simples e rápido de aplicar em relação à amostragem aleatória, mantendo a representatividade da população, especialmente em populações homogêneas.

Amostragem por Conveniência

A **amostragem por conveniência** é o método no qual a amostra é selecionada com base na facilidade de acesso ou disponibilidade dos indivíduos. Esse método é frequentemente utilizado em estudos exploratórios ou quando o objetivo é obter resultados rápidos, mas pode introduzir viés, já que nem todos os membros da população têm a mesma chance de ser selecionados e representados.

Exemplos do Cotidiano



Saúde

Uma clínica pode realizar uma pesquisa de satisfação entre os pacientes que visitam a clínica durante uma semana específica, escolhendo aqueles que estão mais disponíveis.



Economia

Para uma pesquisa de opinião sobre um novo produto, uma empresa pode entrevistar apenas os clientes que visitam a loja em horários de maior movimento.



Agricultura

Um agricultor pode avaliar a produtividade de suas culturas observando apenas as plantas que estão próximas à entrada da fazenda, pois são mais fáceis de acessar.

Esse método é rápido e barato, mas pode não ser representativo da população como um todo, o que limita a validade das conclusões.

Cada um desses métodos de amostragem oferece vantagens e desvantagens, dependendo do objetivo da pesquisa e dos recursos disponíveis. Escolher o método correto é fundamental para garantir a validade dos resultados e a representatividade da amostra em relação à população total.

Para facilitar nosso entendimento sobre os métodos de amostragem, veja na tabela a seguir uma comparação entre eles.

Método de Amostragem	Definição	Como Funciona	Vantagens	Desvantagens
Aleatória Simples	Cada membro da população tem a mesma chance de ser selecionado.	Seleciona-se os indivíduos aleatoriamente, sem padrões.	Sem viés, todos têm a mesma chance de ser escolhidos.	Pode ser difícil de aplicar em grandes populações.
Estratificada	A população é dividida em subgrupos (estratos) e uma amostra é retirada de cada um.	Divide-se a população em estratos e seleciona-se uma amostra proporcional de cada grupo.	Representa bem subgrupos da população, garantindo diversidade.	Exige conhecimento prévio dos estratos e maior esforço na divisão.

Sistemática	Seleciona indivíduos em intervalos regulares após um ponto inicial aleatório.	Seleciona-se o primeiro membro aleatoriamente e depois em intervalos fixos.	Simple e rápido de aplicar, especialmente em populações homogêneas.	Pode introduzir viés se houver um padrão na população que coincida com o intervalo.
Por Conveniência	A amostra é selecionada com base na facilidade de acesso ou disponibilidade.	Seleciona-se os indivíduos que estão mais disponíveis.	Rápido e barato de aplicar.	Pode introduzir viés e não ser representativo da população.

Depois de conhecer os diferentes métodos de amostragem você poderá selecionar o método mais adequado para diferentes contextos e para diferentes tipos de estudos.

Média da População vs. Média da Amostra

Você viu no início desse módulo que a média é uma das ferramentas estatísticas mais utilizadas para resumir a tendência central de um conjunto de dados. Agora vamos ver esta medida usada na população e na amostra.

Quando lidamos com dados, podemos calcular a média de uma população inteira ou de uma amostra, que é um subconjunto dessa população. Compreender a diferença entre a média da população e a média da amostra é essencial para ciência de dados.

Frequentemente, não temos acesso a todos os dados de uma população completa, então trabalhamos com amostras. A partir das médias dessas amostras, podemos fazer estimativas ou inferências sobre o comportamento da população como um todo. Esse conceito é amplamente aplicável em áreas como saúde pública, economia, agricultura e muitos outros campos em que a coleta de dados de toda a população seria impraticável.

Média da População

A média da população representa a média de todos os indivíduos ou elementos que compõem a população inteira. Quando temos acesso a todos os dados da população, podemos calcular uma média precisa, chamada de média populacional. Esta medida reflete o comportamento real da população, sendo o valor verdadeiro em torno do qual os dados da população estão distribuídos.

A fórmula para a média da população é:

$$\mu = \frac{\sum_{i=1}^N x_i}{N}$$

Onde:

μ = média populacional,

x_i = cada valor na população,

N = número total de elementos na população.

Se estamos calculando a média de altura de uma população de 1000 pessoas e temos a altura de cada indivíduo, usamos essa fórmula para somar todas as alturas e dividir pelo número total de pessoas (1000). A média obtida será a média exata da população. Entretanto, em muitos casos práticos, coletar dados de toda uma população é inviável. Seja pelo custo, tempo ou dificuldade de acessar todos os indivíduos, torna-se mais comum o uso de amostras para fazer estimativas.

Veja a seguir exemplos de aplicação da média da população em diferentes áreas.

Exemplos de Média da População



Saúde

Um estudo que mede o peso médio de todos os pacientes atendidos em um hospital durante o ano. Se todos os pacientes forem incluídos, a média obtida será a média da população.



Economia

Ao calcular o salário médio de todos os trabalhadores de uma empresa, considerando todos os funcionários, estamos lidando com a média da população.



Agricultura

O cálculo da altura média de todas as plantas de uma fazenda. Se medirmos a altura de todas as plantas, teremos a média populacional.

Média da Amostra

A média da amostra é a média de um subconjunto da população. Em vez de coletar dados de todos os elementos de uma população, trabalhamos com uma amostra representativa e, a partir dessa amostra, calculamos uma média para estimar a média populacional. A média da amostra é amplamente utilizada em estudos estatísticos e científicos, pois permite que façamos inferências sobre uma população maior com base em dados limitados.

A fórmula para a média da amostra é:

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

Onde:

\bar{x} = média populacional,

x_i = cada valor na população,

n = número total de elementos na população.

Suponha que estamos estudando o peso médio de uma população de 10.000 pessoas, mas em vez de medir todas elas, escolhemos uma amostra aleatória de 200 pessoas. A média do peso dessas 200 pessoas nos fornecerá uma estimativa da média do peso da população. Embora essa média não seja a média exata da população, ela nos dá uma aproximação útil. Veja alguns exemplos práticos a seguir.

Exemplos de Média da Amostra



Saúde

Em vez de medir o peso de todos os pacientes, uma amostra de 100 pacientes pode ser usada para estimar o peso médio de todos os pacientes atendidos no hospital.



Economia

Uma pesquisa pode entrevistar 200 trabalhadores de uma cidade para estimar o salário médio da população total de trabalhadores da cidade.



Agricultura

Se quisermos estimar a altura média das plantas de soja de uma fazenda, podemos medir a altura de 50 plantas em diferentes áreas da plantação.

Podemos perceber que a principal diferença entre a média da população e a média da amostra está no fato de que a média da amostra é apenas uma estimativa da média da população. Quanto mais representativa for a amostra, maior a probabilidade de a média da amostra se aproximar da média da população. Isso é especialmente importante em estudos onde não podemos acessar a população completa, como em pesquisas com grandes populações ou em situações em que há restrições de tempo e custo.

No contexto da ciência de dados, essa prática é fundamental, permitindo que tomemos decisões e façamos previsões com base em amostras, sem a necessidade de dados completos.

Entender e aplicar corretamente esses conceitos, nos ajuda a lidar com incertezas e extrair insights valiosos a partir de conjuntos de dados limitados.

Variáveis e seus Tipos

Uma **variável** é qualquer característica que pode assumir diferentes valores. Por exemplo, altura, idade, cor dos olhos e nível de escolaridade são variáveis porque podem variar entre indivíduos.

As **variáveis** são fundamentais na estatística e na ciência de dados, porque representam qualquer característica, número ou quantidade que pode ser medida ou contada. Identificar corretamente o tipo de variável em uma pesquisa ou estudo é fundamental, pois o tipo de dado determina a escolha do método de análise e interpretação.

Os dois principais tipos de variáveis são as **qualitativas (categóricas)** e as **quantitativas (numéricas)**, além de suas subdivisões.

Variáveis Qualitativas (Categóricas)

As **variáveis qualitativas**, também conhecidas como **categóricas**, descrevem qualidades ou categorias. Elas não podem ser medidas numericamente, mas são classificadas em diferentes grupos ou categorias.

Existem dois tipos de variáveis qualitativas:

- **Nominal:** Quando as categorias não possuem uma ordem lógica ou hierarquia.
- **Ordinal:** Quando as categorias têm uma ordem ou hierarquia.

Para entender melhor, veja alguns exemplos mais comuns em diferentes áreas.

Área	Tipo	Exemplo
Saúde	Nominal	O tipo de sangue (A, B, AB, O) é uma variável qualitativa nominal.
	Ordinal	O nível de dor (leve, moderada, severa) é uma variável qualitativa ordinal.
Economia	Nominal	O setor de atuação de empresas (varejo, manufatura, serviços) é nominal.
	Ordinal	A classificação de crédito (ruim, regular, bom, excelente) é ordinal.
Agricultura	Nominal	As variedades de planta (milho, trigo, soja) são uma variável nominal.
	Ordinal	A qualidade do solo (pobre, média, boa, excelente) é uma variável ordinal.

Variáveis Quantitativas (Numéricas)

As **variáveis quantitativas** são representadas por números e expressam quantidades. Essas variáveis podem ser divididas em dois tipos:

- **Discreta:** Variáveis que assumem valores contáveis e inteiros.
- **Contínua:** Variáveis que podem assumir qualquer valor em um intervalo, incluindo frações e decimais.

Veja alguns exemplos dessas variáveis:

Área	Tipo	Exemplo
Saúde	Discreta	O número de consultas médicas realizadas por um paciente no último ano é uma variável quantitativa discreta, pois o número de consultas é contável e inteiro (ex.: 3 consultas).
	Contínua	O peso de um paciente é uma variável quantitativa contínua, pois pode ser medido com alta precisão e assumir valores fracionários (ex.: 68,5 kg).
Economia	Discreta	O número de produtos vendidos por uma loja em um dia é uma variável quantitativa discreta, pois o número de vendas é contável (ex.: 100 produtos).
	Contínua	O lucro mensal de uma empresa é uma variável quantitativa contínua, pois pode assumir valores com decimais (ex.: R\$ 25.678,75).
Agricultura	Discreta	O número de colheitas em um ano é uma variável quantitativa discreta, pois é contável (ex.: 2 colheitas).
	Contínua	A quantidade de chuva em milímetros registrada em uma fazenda é uma variável quantitativa contínua, pois pode ser medida com precisão fracionada (ex.: 120,8 mm).

A identificação correta do tipo de variável é fundamental para determinar quais métodos estatísticos devem ser usados. Por exemplo:

- Para variáveis qualitativas, podem ser usados gráficos de barras e tabelas de frequências.
- Para variáveis quantitativas, são utilizados gráficos como histogramas e medidas como média e desvio padrão.

Exemplo Prático: Em uma pesquisa agrícola sobre o rendimento de diferentes tipos de solo, um pesquisador poderia coletar dados sobre:

- **Variáveis qualitativas:** O tipo de solo (argiloso, arenoso, misto).
- **Variáveis quantitativas:** A quantidade de grãos colhidos em quilogramas (kg) por hectare, que seria uma variável quantitativa contínua.

Compreender a natureza das variáveis é o primeiro passo para a realização de uma análise de dados eficaz. Dependendo do tipo de variável, diferentes ferramentas e técnicas estatísticas são aplicadas para extrair insights e tomar decisões informadas.

Tipo de Variável	Subtipo	Definição	Exemplo
Qualitativa (Categórica)	Nominal	Categorias sem ordem específica.	Tipo de sangue (A, B, AB, O)
	Ordinal	Categorias com ordem definida.	Nível de dor (leve, moderada, severa)
Quantitativa (Numérica)	Discreta	Valores contáveis e inteiros.	Número de consultas médicas no último ano
	Contínua	Valores em um intervalo, incluindo frações.	Peso de um paciente (68,5 kg)

Além das variáveis qualitativas e quantitativas que já discutimos, existem outros tipos de variáveis que são importantes em estatística e análise de dados.

Veja na tabela a seguir, alguns exemplos de outros tipos de variáveis.

Tipo de Variável	Definição	Exemplo
Binária (Dicotômica)	Duas categorias ou valores possíveis.	Sexo (masculino/feminino), resultado de um teste (aprovado/reprovado)
Multinomial	Mais de duas categorias.	Cor do carro (vermelho, azul, verde, etc.)
Intervalo	Diferenças significativas sem ponto zero absoluto.	Temperatura em graus Celsius
Razão	Diferenças significativas com ponto zero absoluto.	Peso, altura, renda
Contínuas	Qualquer valor em um intervalo.	Altura, temperatura

Discretas	Valores contáveis e inteiros.	Número de filhos, consultas médicas
Tempo	Medida de duração ou momento de eventos.	Tempo gasto em uma tarefa, idade de uma pessoa

Distribuição de Frequência

A **distribuição de frequência** é uma forma de organizar dados para mostrar quantas vezes cada valor ou intervalo de valores ocorre em um conjunto de dados. Ela ajuda a identificar padrões e tendências, facilitando a interpretação dos dados e a visualização de como eles estão distribuídos.

A distribuição de frequência funciona da seguinte forma:

- ▶ Os dados são agrupados em **categorias** (ou intervalos) e o número de vezes que cada categoria ocorre é contado. A contagem de ocorrências em cada categoria é chamada de **frequência**.

Distribuição de frequência simples

Apresenta os valores únicos e a quantidade de vezes que eles aparecem.

Exemplo: Frequência Simples

Imagine que você registrou as idades de 10 pessoas em um grupo:

20, 22, 20, 23, 22, 21, 23, 20, 22, 21

Criando uma distribuição de frequência:

Idade	Frequência
20	3
21	2
22	3
23	2

Os dados na tabela mostram que a idade mais comum no grupo é **20 e 22 anos**, ambas ocorrendo 3 vezes.

Distribuição de frequência agrupada

Agrupa os dados em intervalos quando há muitos valores diferentes, facilitando a análise.

Exemplo: Frequência Agrupada

Agora imagine que você tem os pesos (em kg) de 50 pessoas e os valores variam de 50 a 100. Para facilitar, você pode agrupar os dados em intervalos:

Intervalo de Peso (kg)	Frequência
50-59	8
60-69	15
70-79	12
80-89	10
90-99	5

Esses dados mostram que a maior parte das pessoas tem peso entre **60 e 69 kg**, com 15 ocorrências.

Veja como a distribuição de frequência pode ser usada:

Na exploração dos dados

Antes de usar técnicas avançadas, a distribuição de frequência permite entender a dispersão e os padrões iniciais.

Exemplo: Para analisar quantos clientes compraram produtos em diferentes faixas de preço.

Na detecção de outliers

Valores que aparecem com frequência muito baixa podem indicar anomalias.

Exemplo: Um sistema detecta poucos cliques em uma página web específica, sugerindo problemas.

Fórmulas para calcular frequências

A fórmula para calcular a **frequência relativa** ou **frequência percentual** de um valor em uma **distribuição de frequência** é:

$$\text{Frequência Relativa} = \frac{\text{Frequência Absoluta do Valor}}{\text{Total de Valores}}$$

Se desejar em **percentual**, multiplica-se por 100:

$$\text{Frequência Percentual} = \left(\frac{\text{Frequência Absoluta do Valor}}{\text{Total de Valores}} \right) \times 100$$

Essas fórmulas são para:

- ▶ **Comparar categorias:** É usada para identificar a proporção de um valor ou intervalo específico em relação ao total, permitindo comparações entre categorias.
Exemplo: Saber o percentual de clientes que compraram produtos de uma faixa de preço específica em uma loja.
- ▶ **Fazer análise de dados agrupados:** É aplicada em tabelas de frequência para entender melhor a distribuição e destacar categorias dominantes.
- ▶ **Realizar análises visuais:** É útil na construção de gráficos como pizza (pie charts) e barras, onde as porcentagens ajudam na análise visual.

Veja um exemplo Prático:

Imagine que você registrou as idades de um grupo de 10 pessoas:

20, 22, 20, 23, 22, 21, 23, 20, 22, 21

Já sabemos a distribuição de frequência simples:

Idade	Frequência Absoluta	Frequência Relativa (%)
20	3	$(3/10) \times 100 = 30\%$
21	2	$(2/10) \times 100 = 20\%$
22	3	$(3/10) \times 100 = 30\%$
23	2	$(2/10) \times 100 = 20\%$

Analisando a tabela interpretamos que:

- 30% do grupo tem 20 ou 22 anos.
- Apenas 20% tem 21 ou 23 anos.

Essa fórmula é usada para calcular a proporção de ocorrências de uma categoria em relação ao total, sendo essencial para comparar categorias de forma padronizada e extrair insights mais claros sobre a distribuição dos dados.

A distribuição de frequência organiza dados brutos em categorias e suas respectivas frequências, facilitando a análise e a interpretação. É uma ferramenta essencial na Ciência de Dados para transformar números em informações úteis e aplicáveis.

Histogramas

Os **histogramas** são uma das representações gráficas mais eficazes para visualizar a distribuição de dados, utilizados para identificar padrões, tendências e a dispersão de dados. Em diversas áreas da ciência, como estatística, economia, biologia e ciência de dados, os histogramas são essenciais para revelar a frequência com que diferentes intervalos de valores aparecem em um conjunto de dados.

Um histograma é um gráfico que representa a distribuição de um conjunto de dados usando barras. Cada barra representa um intervalo (ou classe) de valores, e a altura da barra mostra a frequência ou número de ocorrências dos dados dentro desse intervalo.

! ATENÇÃO!

Ao contrário de gráficos de barras comuns, os histogramas têm barras contíguas, sem espaços entre elas, representando a continuidade dos dados em seus intervalos.

Os histogramas são úteis para entender a forma da distribuição (simétrica, enviesada, uniforme etc.) para identificar *outliers* e valores atípicos e visualizar a dispersão e a concentração de dados em diferentes intervalos.

Veja a seguir um passo a passo para construir um Histograma:

1.

O primeiro passo é organizar o conjunto de dados. Por exemplo, imagine que temos os seguintes dados de alturas de uma amostra de 30 pessoas: 160 cm, 162 cm, 165 cm, 170 cm, 175 cm, etc. Para criar um histograma, é necessário dividir os dados em intervalos ou faixas de valores. No caso das alturas, por exemplo, podemos criar intervalos de 5 cm (160-164 cm, 165-169 cm, 170-174 cm, etc.).

2.

Em seguida, contamos quantos valores do conjunto de dados caem em cada intervalo. Para alturas, podemos ter algo como:

160-164 cm (5 pessoas)		165-169 cm (10 pessoas)		170-174 cm (8 pessoas)		175-179 cm (7 pessoas)
---------------------------	--	----------------------------	--	---------------------------	--	---------------------------

3.

No eixo horizontal do histograma (eixo X), colocamos os intervalos, e no eixo vertical (eixo Y), representamos a frequência de dados para cada intervalo. As barras são desenhadas para cada intervalo com altura correspondente à frequência.

Exemplo:

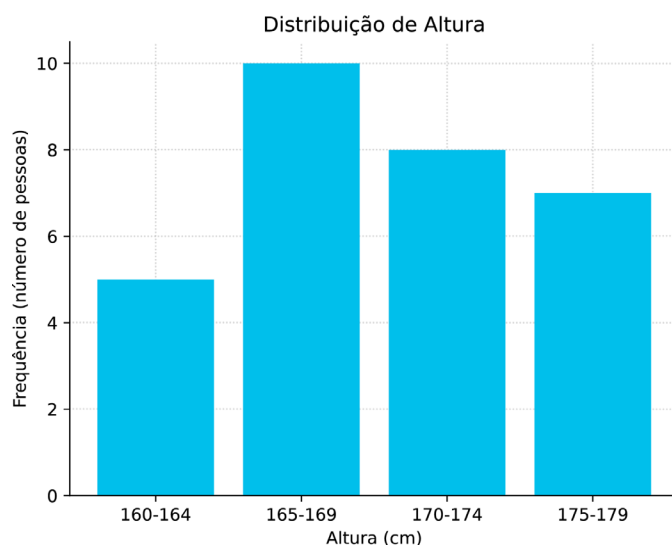
Imagine um conjunto de dados sobre a altura de um grupo de pessoas:

Altura (cm)	Frequência (número de pessoas)
160-164	5
165-169	10
170-174	8
175-179	7

A partir desses dados, podemos desenhar um histograma. Cada intervalo de altura é representado por uma barra, e a altura da barra indica o número de pessoas cujas alturas se enquadram naquele intervalo específico.

Esse histograma nos ajuda a ver imediatamente que a maioria das pessoas neste grupo tem entre 165 e 169 cm de altura.

Figura 6 - Distribuição das alturas por intervalo (cm)



Fonte: Elaboração própria, 2026.

Saber interpretar um histograma é tão importante quanto saber construí-lo. Alguns pontos a serem observados ao analisar histogramas incluem:

- ▶ **Forma da distribuição:** A forma geral da distribuição dos dados pode ser simétrica (como uma curva em forma de sino), assimétrica (com cauda à direita ou esquerda) ou uniforme (quando todas as barras têm altura semelhante).
- ▶ **Outliers:** São valores que se destacam do restante dos dados, localizados em intervalos com baixas frequências.
- ▶ **Dispersão:** A largura das barras e sua distribuição indicam o grau de dispersão dos dados.

MÓDULO 3:

Interpretando Dados com Estatística



Introdução à Interpretação de Dados Estatísticos

Neste módulo, exploraremos conceitos estatísticos fundamentais que ajudam a compreender a distribuição dos dados, como percentis, quantis e o Resumo de Cinco Números. Também abordaremos ferramentas analíticas essenciais, como o **Intervalo Interquartil (IQR)** para detecção de *outliers* e a **Função Densidade de Probabilidade (PDF)**, que auxilia na visualização da dispersão dos valores.

Percentis e Quantis

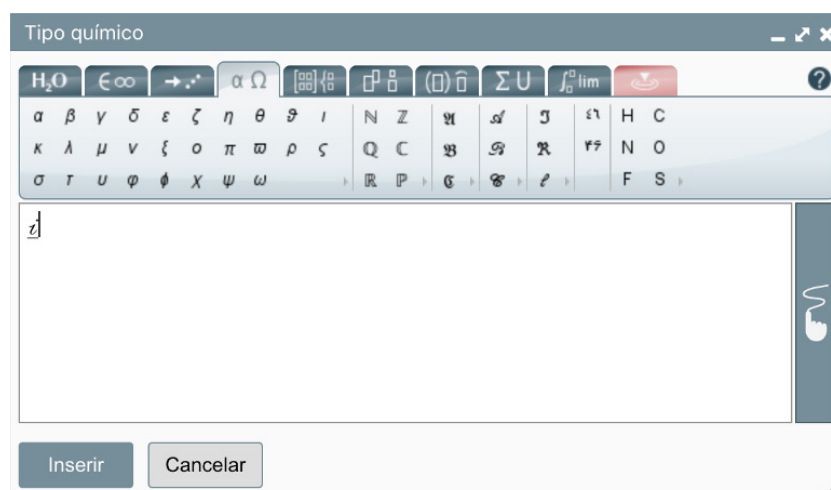
Quando falamos de dados, muitas vezes queremos entender como eles se distribuem e como se comparam entre si. É aqui que entram os percentis e quantis, que são ferramentas muito úteis na estatística. Eles nos ajudam a identificar a posição de um determinado valor em relação aos demais, permitindo que vejamos se esse valor é maior ou menor do que a maioria. Imagine que você fez uma prova e quer saber como sua nota se compara com a de seus colegas. Percentis e quantis podem te dar essa resposta!

Os percentis dividem um conjunto de dados em 100 partes iguais. Por exemplo, se você ouvir que um aluno ficou no percentil 25 (P25) em uma prova, isso significa que 25% dos alunos tiveram notas abaixo dele. Já o percentil 50 (P50) é a mediana, que indica que 50% dos alunos tiveram notas abaixo e 50% acima desse valor. Portanto, saber em que percentil você se encontra pode ser uma boa maneira de entender seu desempenho em relação aos outros.

Por outro lado, os quantis são uma forma mais ampla de dividir os dados. Enquanto os percentis dividem em 100 partes, os quantis podem dividir em qualquer número de partes. Por exemplo, os quartis dividem os dados em quatro partes. Assim, temos o primeiro quartil (Q1), o segundo quartil (Q2, que é a mediana) e o terceiro quartil (Q3).

Veja a seguir como podemos fazer esse cálculo.

Figura 7 - Editor de fórmulas MathType



Fonte: WIRIS, 2025.

Cálculo de Percentis

Calcular um percentil é simples! Veja como você pode fazer isso:

1. Organize os dados: Primeiro, coloque os dados em ordem crescente.
2. Determine o índice do percentil: Use a fórmula:

$$i = \frac{k}{100} \times (n + 1)$$

Onde:

k é o número do percentil desejado,
 n é o número de observações.

Interprete o índice:

- Se i for um número inteiro, o percentil é o valor correspondente ao índice i .
- Se i for um número fracionário, faça uma interpolação linear entre os dois valores adjacentes.

Vamos ver um exemplo prático para ficar mais claro:

Suponha que temos as seguintes notas de um teste:

{ 3,5,7,8,12,14,18,21,22 }

1. Os dados já estão organizados.
2. Vamos calcular o **percentil 40**:

$$i = \frac{40}{100} \times (9 + 1) = 4$$

O índice 4 corresponde ao quarto valor da lista, que é 8. Portanto, o percentil 40 é 8, ou seja, 40% dos alunos obtiveram notas abaixo de 8.

Agora que você já viu como calcular o percentil, vamos entender como calcular o quartil.

Uso de Tabelas para Quantis

Os quantis são frequentemente apresentados em tabelas. A tabela a seguir mostra um exemplo de como percentis e quartis podem ser organizados para um conjunto de dados:

Percentil	Quartil	Valor do Quartil
P25	Q1	5
P50	Q2	8
P75	Q3	18

Essa tabela destaca os quartis, que são os quantis mais comuns usados para dividir os dados em quatro partes.

Aplicações dos Percentis e Quantis

Os percentis e quantis têm várias aplicações práticas:

Na Medicina

Os médicos usam percentis para interpretar resultados de exames. Por exemplo, se seu IMC está no percentil 75, isso significa que você está acima de 75% da população em termos de peso e altura.

Na Educação

Vamos imaginar que você fez um teste padronizado e sua pontuação foi de 670, que corresponde ao percentil 90. Isso significa que você se saiu melhor do que 90% dos outros alunos que fizeram o mesmo teste. É uma ótima maneira de entender seu desempenho!

Em síntese, percentis e quantis são ferramentas poderosas que nos ajudam a entender melhor nossos dados e como eles se comparam. Usá-los pode ser muito útil em diversas áreas, como saúde e educação, ajudando a tomar decisões informadas baseadas em dados concretos. Portanto, da próxima vez que você ouvir sobre percentis em um exame ou teste, saberá exatamente o que isso significa!

Resumo de Cinco Números

O Resumo de Cinco Números é um método importante para descrever um conjunto de dados de maneira simples e eficiente. Essa técnica se baseia em cinco valores que ajudam a entender como os dados estão distribuídos, oferecendo informações sobre a dispersão e a centralidade. Esses valores são: o mínimo, o primeiro quartil, a mediana, o terceiro quartil e o máximo.

1. Mínimo:

O menor valor encontrado no conjunto de dados.

2. Primeiro quartil (Q1):

O valor que separa os 25% menores dados do restante. Também chamado de **percentil 25**.

3. Mediana (Q2):

O valor central que divide o conjunto de dados em duas partes iguais, com 50% dos dados abaixo e 50% acima. Também chamado de **percentil 50**.

4. Terceiro quartil (Q3):

O valor que separa os 75% menores dados dos 25% maiores. Também chamado de **percentil 75**.

5. Máximo:

O maior valor encontrado no conjunto de dados.

Esses cinco números fornecem uma visão clara da distribuição dos dados e são úteis para entender melhor o comportamento deles.

O Resumo de Cinco Números é uma ferramenta poderosa na análise exploratória de dados (AED). Ele permite uma compreensão rápida da distribuição dos dados, proporcionando uma visão geral instantânea. Além disso, essa técnica ajuda a identificar *outliers*, que são valores muito distantes da maioria dos dados. Esses *outliers* podem indicar erros ou fenômenos interessantes.

O resumo também facilita a comparação entre diferentes conjuntos de dados, permitindo analisar suas características de uma só vez. Além disso, ele serve como base para visualizações gráficas, como o *boxplot*, que destaca a distribuição e variação dos dados.

Esse resumo é amplamente utilizado em análises estatísticas e é uma ferramenta poderosa para explorar dados. Ele permite uma compreensão rápida da distribuição e ajuda a identificar pontos que estão muito distantes do resto dos dados, chamados de *outliers*. Esses *outliers* podem indicar erros ou trazer informações importantes que merecem atenção. Além disso, o Resumo de Cinco Números facilita a comparação entre diferentes conjuntos de dados, tornando mais fácil identificar diferenças e semelhanças entre eles.

Outro aspecto interessante é que o Resumo de Cinco Números pode ser visualizado através de um gráfico chamado *boxplot*, que destaca os quartis e os valores mínimos e máximos de forma clara, ajudando a visualizar a dispersão dos dados e a identificar a presença de *outliers* de forma rápida e eficaz.

Na análise exploratória de dados, o Resumo de Cinco Números é frequentemente usado para descrever a distribuição dos dados, fornecendo informações sobre a simetria e dispersão. Ele também prepara o terreno para análises posteriores, ajudando a identificar se os dados precisam de transformações ou ajustes para serem adequadamente analisados por meio de outras técnicas estatísticas.

Vamos considerar um exemplo simples para entender melhor como o Resumo de Cinco Números funciona:

Imagine que temos um conjunto de dados que representa as idades de 15 alunos em uma sala de aula. Esses dados são: 18, 20, 22, 21, 19, 23, 24, 20, 25, 19, 18, 22, 21, 24 e 26. O primeiro passo é organizar esses números em ordem crescente: 18, 18, 19, 19, 20, 20, 21, 21, 22, 22, 23, 24, 24, 25 e 26.

Com os dados organizados, podemos calcular os cinco números que compõem o resumo. O mínimo é 18 e o máximo é 26. A mediana, que é o valor central, será o oitavo número na sequência, ou seja, 21. O primeiro quartil, ou Q1, é a mediana dos primeiros sete números, que é 19. Já o terceiro quartil, ou Q3, é a mediana dos últimos sete números, que é 24. Portanto, o Resumo de Cinco Números para esse conjunto de dados é: Mínimo: 18, Q1: 19, Mediana: 21, Q3: 24 e Máximo: 26.

Essa análise simples já nos dá uma boa ideia de como as idades dos alunos estão distribuídas. Podemos ver, por exemplo, que metade dos alunos tem entre 19 e 24 anos, e que as idades mais extremas são 18 e 26 anos. Se quisermos ir além, podemos usar um *boxplot* para visualizar esses números de maneira gráfica, o que facilita ainda mais a interpretação dos dados.

O Resumo de Cinco Números é uma ferramenta fundamental para qualquer pessoa que queira entender rapidamente um conjunto de dados. Ele ajuda a resumir informações complexas de forma acessível e direta, permitindo que possamos tirar conclusões importantes e tomar decisões com base em uma análise clara e objetiva dos dados.

Intervalo Interquartil (IQR)

O Intervalo Interquartil (ou IQR, em inglês) é uma medida de dispersão que nos ajuda a entender o quanto os dados estão espalhados, focando na parte central do conjunto de dados. Ele é uma ferramenta muito útil porque exclui os valores mais extremos e concentra-se nos dados que estão mais próximos da mediana.

Para calcular o IQR, usamos dois valores importantes: o primeiro quartil (Q1) e o terceiro quartil (Q3). O Q1 é o valor que separa os 25% menores dados, enquanto o Q3 separa os 25% maiores. A fórmula do IQR é simples:

$$IQR = Q3 - Q1$$

Essa fórmula nos mostra a diferença entre o valor no terceiro quartil e o valor no primeiro quartil, ou seja, ela nos diz em que faixa os 50% centrais dos dados estão distribuídos.

Veja um exemplo prático:

Vamos usar um exemplo simples para entender melhor como calcular o IQR. Imagine o seguinte conjunto de dados que representa as notas de 10 alunos:

$$\{12, 15, 14, 10, 20, 18, 25, 30, 27, 22\}$$

O primeiro passo é organizar esses números em ordem crescente:

$$\{10, 12, 14, 15, 18, 20, 22, 25, 27, 30\}$$

Agora, podemos calcular os quartis.

- **Mediana (Q2):** Como temos 10 dados, a mediana será a média entre o 5º e o 6º valores, que são 18 e 20. Assim, a mediana será:

$$Q2 = \frac{18 + 20}{2} = 19$$

- **Primeiro Quartil (Q1):** O Q1 é a mediana dos primeiros cinco valores {10, 12, 14, 15, 18}. Nesse caso, o valor central é 14, então $Q1 = 14$.
- **Terceiro Quartil (Q3):** O Q3 é a mediana dos últimos cinco valores {20, 22, 25, 27, 30}. O valor central aqui é 25, então $Q3 = 25$.

Agora que temos Q1 e Q3, podemos calcular o IQR:

$$IQR = 25 - 14 = 11$$

Isso significa que os 50% centrais dos dados estão distribuídos em uma faixa de 11 unidades. Ou seja, a maioria das notas está concentrada em um intervalo de 11 pontos.

Mas o que significa e o que fazer com os dados fora desse limite? Eles também precisam ser tratados, são os *outliers* e o IQR.

Identificação de *Outliers*

Uma das grandes utilidades do IQR é que ele nos ajuda a identificar *outliers*, que são valores muito diferentes do resto dos dados. Para descobrir se existem *outliers*, calculamos os limites inferior e superior, que nos dizem se algum valor está "fora" da faixa esperada dos dados. Esses limites são calculados da seguinte maneira:

$$\begin{aligned}\text{Limite inferior} &= Q1 - 1.5 \times \text{IQR} \\ \text{Limite Superior} &= Q3 + 1.5 \times \text{IQR}\end{aligned}$$

Se algum valor estiver abaixo do limite inferior ou acima do limite superior, ele será considerado um *outlier*. Vamos aplicar essa ideia ao nosso conjunto de dados:

$$\begin{aligned}\text{IQR} &= 11 \\ \text{Limite inferior} &= 14 - 1.5 \times 11 = -1.5 \\ \text{Limite Superior} &= 25 + 1.5 \times 11 = 40.5\end{aligned}$$

Então, qualquer valor abaixo de -2,5 ou acima de 41,5 seria considerado um *outlier*. No entanto, como as notas do nosso conjunto de dados estão todas entre 10 e 30, podemos concluir que não há *outliers* nesse caso.

Importância do IQR

O IQR é uma medida muito útil porque, ao contrário de outras medidas de dispersão como a variância ou o desvio padrão, ele não é afetado por *outliers*. Isso significa que, se houver valores extremamente altos ou baixos no conjunto de dados, o IQR ainda vai nos dar uma visão precisa da dispersão dos dados centrais. Por isso, o IQR é frequentemente preferido quando queremos entender a distribuição dos dados sem que os extremos distorçam a análise. Além disso, o IQR é a base para a construção de gráficos chamados *boxplots*, que são muito úteis para visualizar a distribuição dos dados e identificar *outliers* de forma clara.

Você não deve esquecer que o Intervalo Interquartil (IQR) é uma ferramenta valiosa para medir a dispersão dos dados e identificar *outliers*. Ele oferece uma visão clara dos dados centrais, sem ser influenciado por valores extremos, e é amplamente utilizado em análises estatísticas. Compreender o IQR é essencial para quem quer analisar dados de forma eficiente e precisa, seja em um contexto escolar ou em pesquisas mais avançadas.

Boxplots

Os *boxplots* são ferramentas visuais simples, mas poderosas, que ajudam a entender a

distribuição dos dados e a identificar possíveis valores extremos, chamados de *outliers*. Eles são amplamente utilizados para resumir dados estatísticos e visualizar a dispersão em torno da mediana de forma clara e eficiente.

Um *boxplot* se baseia no Resumo de Cinco Números (mínimo, Q1, mediana, Q3, e máximo) e utiliza uma caixa e "bigodes" para representar a variação dos dados, enquanto pontos fora do alcance dos bigodes indicam *outliers*.

Elementos de um Boxplot

Vamos recordar os cinco números, pois eles são os principais elementos de um *boxplot*:

- **Mínimo:** O menor valor dos dados que não é um *outlier*.
- **Primeiro Quartil (Q1):** O ponto abaixo do qual 25% dos dados estão localizados.
- **Mediana (Q2):** O valor central que divide os dados em duas partes iguais.
- **Terceiro Quartil (Q3):** O ponto abaixo do qual 75% dos dados estão localizados.
- **Máximo:** O maior valor dos dados que não é um *outlier*.

Os bigodes que saem da caixa do *boxplot* indicam o intervalo entre os quartis e os valores mínimo e máximo. Pontos que se encontram fora deste intervalo são *outliers*.

Seguindo nosso raciocínio, vamos ver como a Dispersão e Identificar Outliers.

Visualizar a Dispersão e Identificar Outliers

Um dos principais benefícios do *boxplot* é sua capacidade de revelar a dispersão dos dados em relação à mediana. A largura da caixa do *boxplot* reflete a variação entre Q1 e Q3, conhecida como Intervalo Interquartil (IQR). Quanto maior a caixa, maior é a dispersão dos dados. Além disso, o *boxplot* é excelente para identificar *outliers*, que são valores muito acima ou abaixo do esperado. Valores que caem fora dos limites dos bigodes são considerados *outliers*.

Esses limites são calculados com base no IQR, usando a seguinte fórmula:

$$\text{Limite inferior} = Q1 - 1.5 \times IQR$$

$$\text{Limite Superior} = Q3 + 1.5 \times IQR$$

Se algum dado estiver fora desses limites, ele é considerado um *outlier*.

Vejamos um exemplo prático:

Vamos usar como exemplo um conjunto de dados que representa as alturas (em cm) de um grupo de estudantes:

Conjunto de dados:

{150, 155, 160, 165, 170, 175, 180, 190, 200, 210}

Passo 1: Organizar os dados em ordem crescente:

{150, 155, 160, 165, 170, 175, 180, 190, 200, 210}

Passo 2: Calcular os quartis e a mediana:

- **Q1 (Primeiro Quartil):** Mediana dos cinco primeiros valores. $Q1=160$.
- **Mediana (Q2):** Mediana de todo o conjunto. Como há 10 valores, a mediana será a média dos 5º e 6º valores. $Q2=(170+175)/2= 175.5$.
- **Q3 (Terceiro Quartil):** Mediana dos últimos cinco valores. $Q3=190$.

Passo 3: Calcular o IQR:

$$IQR = Q3 - Q1$$
$$IQR = 190 - 160 = 30$$

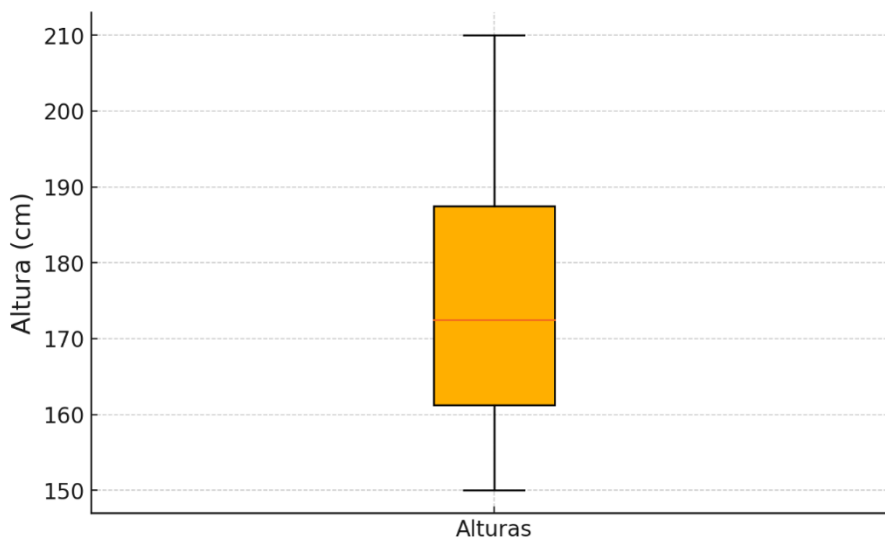
Passo 4: Identificar outliers:

$$\text{Limite inferior} = 160 - 1.5 \times 30 = 115$$

$$\text{Limite Superior} = 190 + 1.5 \times 30 = 235$$

Os valores estão todos dentro desse intervalo, portanto, não temos outliers. Podemos agora visualizar o boxplot baseado nesses dados. A mediana é representada pela linha dentro da caixa, e os bigodes se estendem até os valores mínimo (150) e máximo (210).

Figura 8 - Boxplot das Alturas do Alunos (cm)



Fonte: Elaboração própria, 2025.

O gráfico mostra a distribuição das alturas, destacando a mediana, os quartis e a variação dos dados.

Efeito de Outliers e Sua Remoção

No mundo da estatística, os *outliers*, ou valores extremos, são dados que se afastam significativamente do restante do conjunto.

Figura 9 - Cesta de basquete

Imagine uma situação em que você e seus amigos jogam basquete, e todos marcam entre 5 e 20 pontos por jogo, mas um amigo, em um dia excepcional, marca 50 pontos.

Esse desempenho incomum é um *outlier*.



Fonte: Markus Spiske, pexels, 2026.

Os *outliers* podem ter um grande impacto em nossas análises, distorcendo resultados e levando a conclusões erradas.

Quando calculamos a média de um conjunto de dados, a presença de um *outlier* pode puxar o resultado para cima ou para baixo, tornando a média menos representativa.

Por exemplo, se a média das pontuações de basquete for calculada com o amigo que fez 50 pontos, essa média será muito maior do que a realidade da maioria. Além disso, os *outliers* também afetam a variância, que mede a dispersão dos dados. Com um *outlier*, a variância pode parecer maior do que realmente é, o que dá uma impressão errada sobre a consistência dos dados.

Para detectar *outliers*, uma técnica comum é usar o Intervalo Interquartil (IQR). Já discutimos essa ferramenta anteriormente, e ela pode ser usada para identificar valores que estão muito distantes da maioria. Se um dado estiver a mais de 1,5 vezes o IQR abaixo do primeiro quartil ou acima do terceiro quartil, ele é considerado um *outlier*.

Uma vez identificados, devemos considerar se esses *outliers* devem ser removidos ou não. A remoção pode ser apropriada se os *outliers* forem resultados de erros de medição ou dados que não se encaixam na análise que estamos realizando. Por outro lado, se o *outlier* representa uma variação legítima nos dados, removê-lo pode resultar em perda de informações importantes.

Para ilustrar essa ideia, vamos considerar um exemplo prático:

Suponha que temos as notas finais de uma turma de matemática: {60, 62, 65, 68, 70, 72, 75, 80, 85, 95, 100, 200}. A nota 200 é claramente um *outlier*, pois é muito maior que as demais. Calculando a média e a variância com e sem essa nota, percebemos que a média e a variância mudam drasticamente. Quando removemos o 200, a média das notas passa a ser mais representativa do desempenho real da turma.

Além disso, a remoção do *outlier* pode ser visualizada em gráficos, como *boxplots*, que mostram a distribuição dos dados. Um *boxplot* que inclui o *outlier* terá uma aparência muito diferente de um *boxplot* que o exclui, evidenciando como a presença de um *outlier* pode distorcer nossa percepção dos dados.

Entender os *outliers* e sua remoção é essencial na análise de dados. Eles podem influenciar nossos resultados de maneira significativa, e, portanto, devemos ser cautelosos ao interpretá-los. Identificá-los corretamente e decidir se devem ser mantidos ou removidos é uma habilidade importante para qualquer analista de dados.

Função Densidade de Probabilidade

A função densidade de probabilidade, ou PDF (do inglês, *Probability Density Function*), é um conceito fundamental em estatística que nos ajuda a entender como as variáveis contínuas se distribuem.

Vamos explorar o que é a PDF, como ela se aplica a diferentes distribuições e como interpretar seus gráficos. A PDF é uma função matemática que descreve a probabilidade relativa de uma variável contínua assumir um determinado valor.

Diferentemente de variáveis discretas, onde podemos contar a probabilidade de um resultado específico, em variáveis contínuas lidamos com intervalos. A PDF nos fornece uma maneira de visualizar a probabilidade de que a variável caia dentro de um intervalo específico.

Para entender melhor, pense em uma variável contínua como a altura de pessoas em uma sala. Enquanto uma pessoa pode ter uma altura de exatamente 1,75 metros, existem muitas outras que variam em alturas, e é impossível ter uma medição exata para cada centímetro. A PDF permite que visualizemos a distribuição de alturas, mostrando onde a maioria das alturas se concentra e onde há menos ocorrências.

Vejamos a seguir alguns exemplos mais comuns de distribuição.

Exemplos de Distribuições Comuns

As distribuições de probabilidade mais conhecidas incluem a distribuição normal e a distribuição uniforme.

Distribuição Normal

Essa é uma das distribuições mais importantes em estatística. Ela é representada por uma curva em forma de sino, onde a média, mediana e moda são iguais e estão localizadas no centro. A maioria dos dados está próxima da média, enquanto valores extremos são menos prováveis.

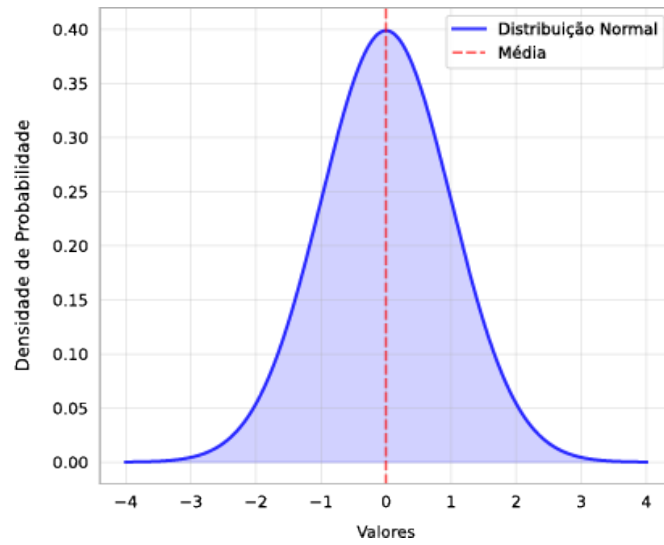
A função densidade de probabilidade da distribuição normal é dada pela fórmula:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Aqui, μ representa a média e σ o desvio padrão da distribuição.

Na figura abaixo, a Distribuição Normal apresenta uma curva simétrica em torno da média (linha vermelha pontilhada). A área sob a curva representa a probabilidade total, que é igual a 1.

Figura 10 - Função de Densidade de Probabilidade - Distribuição Normal



Fonte: Elaboração própria, 2025.

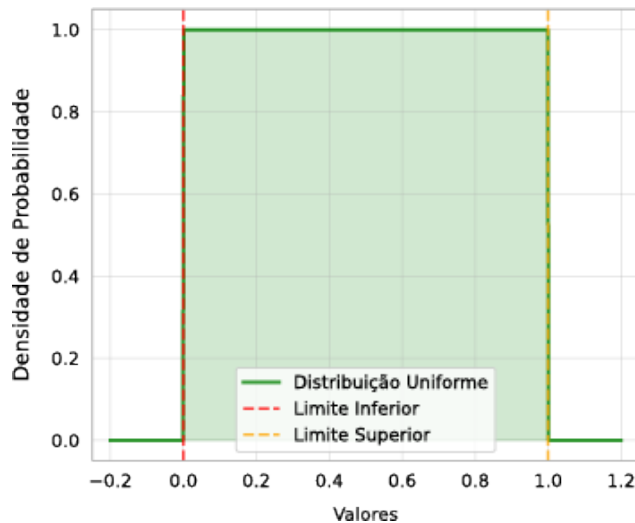
Distribuição Uniforme

Nesta distribuição, todos os resultados têm a mesma probabilidade de ocorrer. Imagine um dado justo: a chance de obter qualquer número de 1 a 6 é a mesma. A função densidade de probabilidade para a distribuição uniforme em um intervalo $[a, b]$ é dada por:

$$f(x) = \frac{1}{b - a} \text{ para } a \leq x \leq b$$

Na figura abaixo, a Distribuição Uniforme representa uma densidade constante entre os limites inferior e superior (linhas vermelhas e laranjas pontilhadas). Todos os valores dentro desse intervalo têm a mesma probabilidade.

Figura 11 - Função de Densidade de Probabilidade - Distribuição Uniforme

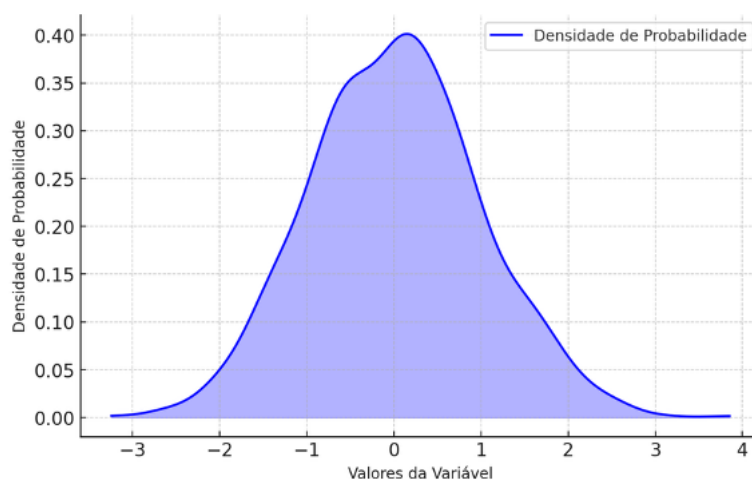


Fonte: Elaboração própria, 2025.

Interpretação de Gráficos de Densidade

Os gráficos de densidade são ferramentas valiosas para visualizar a distribuição de dados contínuos. O **eixo x representa os valores da variável, enquanto o eixo y representa a densidade de probabilidade**. É importante lembrar que a área sob a curva em um intervalo específico representa a probabilidade de a variável cair nesse intervalo. Por exemplo, em uma distribuição normal, a área sob a curva entre um desvio padrão acima e abaixo da média contém aproximadamente 68% dos dados. Ao analisar um gráfico de densidade, podemos observar a forma da distribuição (se é simétrica ou enviesada), a presença de picos (modos) e como a densidade varia. Identificar essas características é essencial para realizar inferências sobre os dados e aplicar modelos estatísticos adequados.

Figura 12 - Gráfico de Densidade de uma Distribuição Normal



Fonte: Elaboração própria, 2025.

- **Eixo X (horizontal):** Representa os **valores da variável**, variando aproximadamente de -3 a 4. Ele indica os possíveis valores que a variável pode assumir dentro da distribuição normal gerada.
- **Eixo Y (vertical):** Representa a **densidade de probabilidade**, variando de 0 a cerca de 0.40. Esse eixo mostra a probabilidade relativa de cada valor no eixo X ocorrer.

A curva azul mostra a distribuição dos dados, com um pico central indicando que a maior parte dos valores está concentrada próximo da média (zero, no caso da distribuição normal gerada).

Podemos perceber que, a função densidade de probabilidade é uma ferramenta poderosa que nos ajuda a entender e visualizar distribuições de variáveis contínuas. Ao conhecer as distribuições comuns, como a normal e a uniforme, e ao interpretar gráficos de densidade, somos capazes de fazer análises mais precisas e tomar decisões informadas em diversas áreas, como ciências sociais, engenharia e medicina. Compreender a PDF é um passo fundamental na jornada para se tornar proficiente em estatística e análise de dados.

Pontuação Z

A pontuação Z é uma medida estatística que indica quantos desvios padrão um determinado valor está acima ou abaixo da média de um conjunto de dados. Ela nos permite comparar valores em diferentes distribuições, mesmo que essas distribuições tenham médias e desvios padrão diferentes.

A fórmula para calcular a pontuação Z é:

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

Onde:

Z é a pontuação Z,

X é o valor observado,

μ é a média do conjunto de dados,

σ é o desvio padrão do conjunto de dados.

Vejamos na prática como fazer o cálculo da Pontuação Z:

Vamos considerar um exemplo prático para entender como calcular a pontuação Z.

Exemplo 1: Suponha que temos as notas de uma prova em uma turma:

Notas: {70, 75, 80, 85, 90}

1º Passo: Calcule a média (μ):

$$\sigma = \frac{70 + 75 + 80 + 85 + 90}{5} = \frac{400}{5} = 80$$

2º Passo: Calcule o desvio padrão (σ):

$$\sigma = \sqrt{\frac{(70 - 80)^2 + (75 - 80)^2 + (80 - 80)^2 + (85 - 80)^2 + (90 - 80)^2}{5}}$$

$$\sigma = \sqrt{\frac{100 + 25 + 0 + 25 + 100}{5}}$$

$$\sigma = \sqrt{\frac{250}{5}}$$

$$\sigma = \sqrt{50}$$

$$\sigma \approx 7,07$$

3º Passo: Calcule a pontuação Z para uma nota de 85:

$$Z = \frac{(85 - 80)}{7,07} \approx \frac{5}{7,07} \approx 0,71$$

Isso significa que a nota 85 está aproximadamente 0,71 desvios padrão acima da média.

Exemplo 2: Agora, vamos considerar uma nota de 70

$$Z = \frac{(70 - 80)}{7,07} \approx \frac{-10}{7,07} \approx -1,41$$

Isso indica que a nota 70 está aproximadamente 1,41 desvios padrão abaixo da média.

Aplicações da Pontuação Z

A pontuação Z tem várias aplicações práticas em estatística e análise de dados:

Detecção de Outliers

Valores que têm uma pontuação Z maior que 3 ou menor que -3 são geralmente considerados outliers. Esses valores podem ser extremos e merecem uma investigação mais aprofundada.

Padronização

A pontuação Z permite a padronização de diferentes conjuntos de dados, facilitando comparações entre eles, mesmo quando as escalas são diferentes. Por exemplo, se quisermos comparar o desempenho de alunos em duas provas com escalas diferentes, podemos usar a pontuação Z para torná-las comparáveis.

A pontuação Z é uma ferramenta valiosa na análise estatística, pois nos ajuda a entender a posição de um valor em relação a uma média e a identificar valores extremos. Compreender como calcular e interpretar a pontuação Z é essencial para qualquer estudante que deseja se aprofundar em estatística e análise de dados.

- ▶ Quando lidamos com dados em ciência de dados, frequentemente encontramos variáveis com escalas diferentes. Isso pode dificultar a análise, especialmente em algoritmos que dependem da distância entre os dados, como a regressão linear e os métodos de agrupamento.
- ▶ Duas técnicas comuns para ajustar os dados a uma escala comparável são a **padronização** e a **normalização**.

Padronização x Normalização

A padronização é uma técnica que ajusta os dados para que tenham uma média de 0 e um desvio padrão de 1. Isso significa que, ao padronizar os dados, transformamos a distribuição de valores para que se concentre em torno de zero, com uma dispersão uniforme.

A fórmula usada para padronizar um valor x é:

$$z = \frac{x - \mu}{\sigma}$$

Onde:

Z o valor padronizado (ou pontuação Z),

x é o valor original,

μ é a média dos dados,

σ é o desvio padrão dos dados.

Quando usar padronização?

A padronização é recomendada quando os dados seguem uma distribuição normal ou quando os algoritmos que você utiliza assumem que os dados possuem essa distribuição. Ela é muito comum em técnicas estatísticas como a regressão linear e na maioria dos métodos de aprendizado de máquina supervisionado.

Vejamos um exemplo prático:

Imagine que temos os seguintes dados de alturas de um grupo de estudantes (em centímetros):

{160, 170, 180, 190, 200}

A média μ das alturas é 180, e o desvio padrão σ é 15. Vamos padronizar o valor 190:

$$z = \frac{190 - 180}{15} = \frac{10}{15} = 0,67$$

Isso significa que a altura de 190 cm está 0.67 desvios-padrão acima da média.

Normalização

A normalização, por outro lado, é uma técnica que ajusta os dados para que eles fiquem em uma escala fixa, geralmente entre 0 e 1. Isso é útil quando os dados têm valores em escalas muito diferentes e você deseja torná-los diretamente comparáveis.

A fórmula usada para normalizar um valor x é:

$$\acute{x} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

Onde:

\acute{x} o valor normalizado,

x é o valor original,

x_{min} é valor mínimo do conjunto de dados,

x_{max} é o valor máximo do conjunto de dados.

Quando usar normalização?

A normalização é mais apropriada quando não sabemos a distribuição dos dados ou quando os algoritmos que vamos utilizar são baseados em distância, como os k-vizinhos mais próximos (k-NN) e redes neurais. Ela é útil quando precisamos de dados em uma escala comparável.

Veja um exemplo de normalização:

Vamos usar o mesmo conjunto de dados de alturas:

{160, 170, 180, 190, 200}.

O valor mínimo é 160 e o valor máximo é 200. Vamos normalizar o valor 190:

$$\hat{x} = \frac{190 - 160}{200 - 160} = \frac{30}{40} = 0,75$$

Isso significa que 190 está 75% do caminho entre o valor mínimo e o máximo.

Tanto a padronização quanto a normalização são técnicas essenciais na preparação dos dados, mas devem ser aplicadas de acordo com a natureza dos dados e os algoritmos que você deseja utilizar. A padronização é mais indicada quando os dados seguem uma distribuição normal, enquanto a normalização é mais útil para escalar dados para comparações diretas. Compreender essas técnicas ajuda a evitar distorções nas análises e melhora a eficácia dos modelos de aprendizado de máquina.

Referências

Han, J., Kamber, M., & Pei, J. (2012). *Data Mining: Concepts and Techniques* (3rd ed.) Elsevier.

Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (2nd ed.) Springer.

Laney, D. (2001). *3D Data Management: Controlling Data Volume, Velocity, and Variety*. META Group Research.

Lima, M. (2023). Medidas de dispersão: A amplitude, a variância e o desvio-padrão. *Blog Psicometria Online*. <https://www.blog.psicometriaonline.com.br/medidas-de-dispersao-amplitude-a-variancia-e-o-desvio-padrao/>. Acesso em: 26 nov. 2024.

McMahan B., Ramage D. (2017) Federated Learning: Collaborative Machine Learning without Centralized Training Data. *Google Research*. <https://research.google/blog/federated-learning-collaborative-machine-learning-without-centralized-training-data/>. Acesso em: 26 nov. 2024.

Microsoft. (2025). O que é AutoML (machine learning automatizado)? *Documentação do Azure Machine Learning*. <https://learn.microsoft.com/pt-br/azure/machine-learning/concept-automated-ml?view=azureml-api-2>. Acesso em: 26 nov. 2024.

Naur, P. (1974). *Concise Survey of Computer Methods*. Academic Press.

Fawcett, T., & Provost, F. (2013). *Data Science for Business: What You Need to Know about Data Mining and Data-Analytic Thinking*. O'Reilly Media.

Cleveland, W. S. (2001). Data Science: An Action Plan for Expanding the Technical Areas of the Field of Statistics. *International Statistical Review*, 69(1), 21-26.

