

# Revisão e Notas sobre Estatística I

Pedro Rupf Pereira Viana

12 de dezembro de 2025

# 1 Introdução e Objetivos

Trata-se de uma revisão teórica e expositiva sobre Conceitos gerais, Análise Exploratória de Dados (EDA), Amostragem, Probabilidade, Distribuição Normal e o Teorema Central do Limite (TCL), bem como suas aplicações na área de Análise e Ciência de Dados.

## 2 Conceitos Gerais

A palavra **Estatística** têm sua origem na palavra latina *STATUS* (Estado). Há indícios de que há 3000 anos A.C. já eram realizados censos na Babilônia, China e Egito. Também podemos ver no livro de Números, do Velho Testamento, uma instrução do Altíssimo dada à Moisés, para que fosse realizado um levantamento dos homens de Israel que estivessem aptos para guerrear (Nm. 1:1-3). Usualmente, estas informações eram utilizadas para a taxação de impostos ou para o alistamento militar.

A estatística envolve técnicas para coletar, organizar, descrever, analisar e interpretar dados, quer sejam provenientes de experimentos, ou vindos de estudos observacionais. A análise estatística de dados geralmente tem por objetivo a tomada de decisões, resoluções de problemas ou produção de conhecimento. Mas novos conhecimentos normalmente geram problemas de pesquisas, resultando em um processo iterativo.

Normalmente o termo estatística esta associado a números, tabelas, gráficos, mas a importância da estatística fica melhor representada por dois pontos comuns em nosso dia a dia:

- Dados
- Variabilidade

Podemos dividir a estatística em três áreas:

- Estatística Descritiva
- Probabilística
- Estatística Inferencial

A estatística descritiva é a etapa inicial da análise utilizada para descrever e resumir os dados. A disponibilidade de uma grande quantidade de dados e de métodos computacionais muito eficientes revigorou esta área da estatística.

Inferência Estatística é, em suma, o estudo de técnicas fundamentada na teoria das probabilidades, que possibilitam a extrapolação a um grande conjunto de dados, das informações e conclusões obtidas a partir da amostra.

### 2.1 Estatística Descritiva

A estatística descritiva é um ramo fundamental da estatística que se concentra em resumir, organizar e apresentar dados de forma clara e concisa, sem fazer inferências sobre populações maiores. Ela permite que pesquisadores, analistas e tomadores de decisão compreendam as características principais de um conjunto de dados, facilitando a identificação de padrões, tendências e anomalias. De acordo com definições padrão, a estatística descritiva utiliza coeficientes breves para sintetizar um conjunto de dados, representando

uma população inteira ou uma amostra dela. Essa abordagem é essencial em diversas áreas, como economia, saúde, educação e ciências sociais, onde o volume de dados pode ser esmagador sem ferramentas de síntese.

Historicamente, a estatística descritiva evoluiu a partir de métodos simples de contagem e agrupamento, ganhando sofisticação com o avanço da computação. Hoje, ela é distinguida da estatística inferencial, que usa amostras para fazer previsões sobre populações. Em resumo, enquanto a inferencial projeta para o desconhecido, a descritiva descreve o conhecido. Para ilustrar, imagine um conjunto de dados sobre as alturas de alunos em uma turma: a descritiva calcularia médias e variações, mas não generalizaria para todos os alunos do país.

### 2.1.1 Tipos de dados e variáveis em Estatística Descritiva

Antes de mergulhar nas medidas, é crucial entender os tipos de dados manipulados na estatística descritiva. Os dados podem ser classificados como qualitativos (categóricos) ou quantitativos (numéricos). Os qualitativos incluem categorias nominais (sem ordem, como cores) e ordinais (com ordem, como níveis de satisfação: baixo, médio, alto). Já os quantitativos dividem-se em discretos (valores inteiros, como número de filhos) e contínuos (valores reais, como peso).

Essa classificação influencia as técnicas descritivas aplicadas. Por exemplo, para dados qualitativos, usamos frequências e percentuais, enquanto para quantitativos, aplicamos medidas numéricas. Além disso, as variáveis podem ser independentes (causadoras) ou dependentes (afetadas), embora na descritiva o foco seja mais na descrição do que na causalidade. Uma análise descritiva eficaz começa com a identificação correta desses tipos, evitando erros na interpretação.

### 2.1.2 Medidas de Tendência Central

As medidas de tendência central são o cerne da estatística descritiva, resumindo um conjunto de dados em um único valor representativo. As principais são a **média**, a **mediana** e a **moda**.

- **Média Aritmética:** Calculada como a soma de todos os valores, divididos pelo número de observações:

$$\bar{x} = \frac{\sum x_i}{n} \quad (2.1)$$

Ela é sensível a valores extremos (chamados de *outliers*), o que a torna útil para dados simétricos, mas menos robusta em distribuições assimétricas. Por exemplo, em um conjunto de salários 1000, 1200, 1300, 50000, a média é alta devido ao outlier.

- **Mediana:** O valor central quando os dados são ordenados. Para um número ímpar de observações, é o meio; para par, a média dos dois centrais. Ela é resistente a outliers, ideal para dados enviesados, como rendas familiares.
- **Moda:** O valor mais frequente. Pode ser unimodal (uma moda), bimodal (duas) ou multimodal. Útil para dados categóricos, como a cor mais comum em uma pesquisa.

Essas medidas fornecem insights sobre o "centro" dos dados, mas devem ser usadas em conjunto para uma visão completa. Em aplicações reais, como análise de desempenho escolar, a mediana pode revelar uma tendência mais realista do que a média influenciada por notas extremas.

### 2.1.3 Medidas de Dispersão

Enquanto as medidas de tendência central indicam onde os dados se concentram, as de dispersão mostram o quão espalhados eles estão. Elas quantificam a variabilidade, essencial para entender a consistência dos dados.

- **Amplitude:** A diferença entre o valor máximo e mínimo ( $\text{Amplitude} = \text{Máx} - \text{Mín}$ ). Simples, mas sensível a outliers.
- **Variância:** A média dos quadrados das diferenças em relação à média:

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1} \quad (2.2)$$

Mede a dispersão quadrática, útil em análises financeiras para risco.

- **Desvio Padrão:** A raiz quadrada da variância ( $s = \sqrt{s^2}$ ). Na mesma unidade dos dados, facilita interpretações, como em controle de qualidade industrial.
- **Intervalo Interquartil (IQR):** A diferença entre o terceiro quartil (Q3) e o primeiro (Q1), capturando a dispersão do meio 50% dos dados. Robusto contra outliers, usado em boxplots.

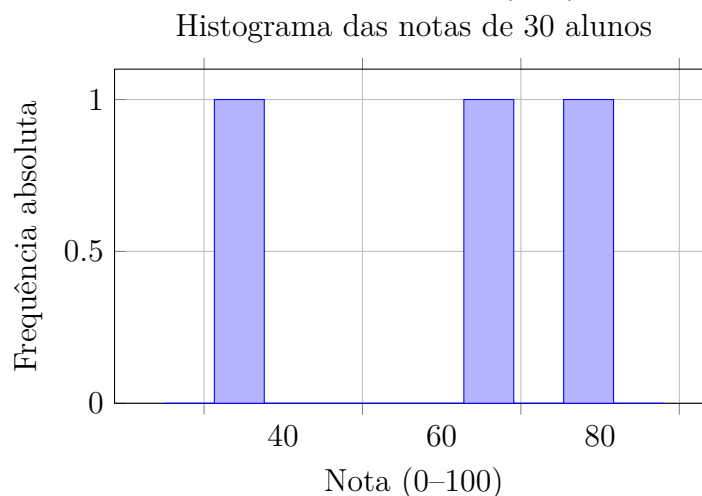
Essas medidas ajudam a comparar conjuntos de dados. Por exemplo, duas turmas com a mesma média de notas podem ter desvios padrões diferentes, indicando variabilidades distintas.

### 2.1.4 Distribuições de Frequência e Representações Gráficas

A estatística descritiva não se limita a números; ela emprega gráficos para visualização. As distribuições de frequência agrupam dados em classes, mostrando contagens ou percentuais. Por exemplo, uma tabela de frequência para idades pode revelar picos em faixas etárias.

Gráficos comuns incluem:

- **Histograma:** Barras representando frequências em intervalos, ideal para dados contínuos, revelando formas como normal (sino) ou assimétrica.



Histograma com 10 classes

- **Gráfico de Barras:** Para dados categóricos, comparando frequências.  
Preferência de sabor de sorvete (n=200)

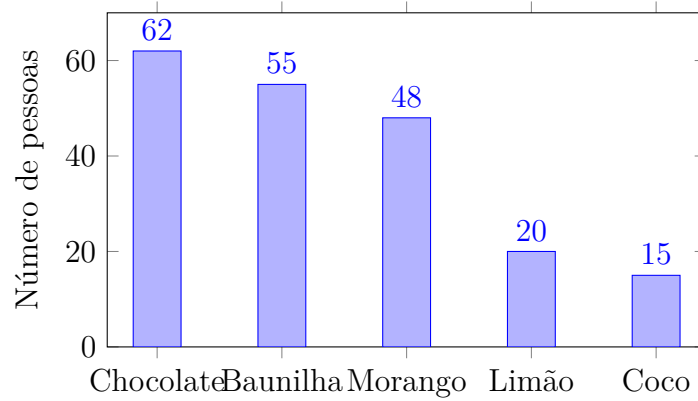
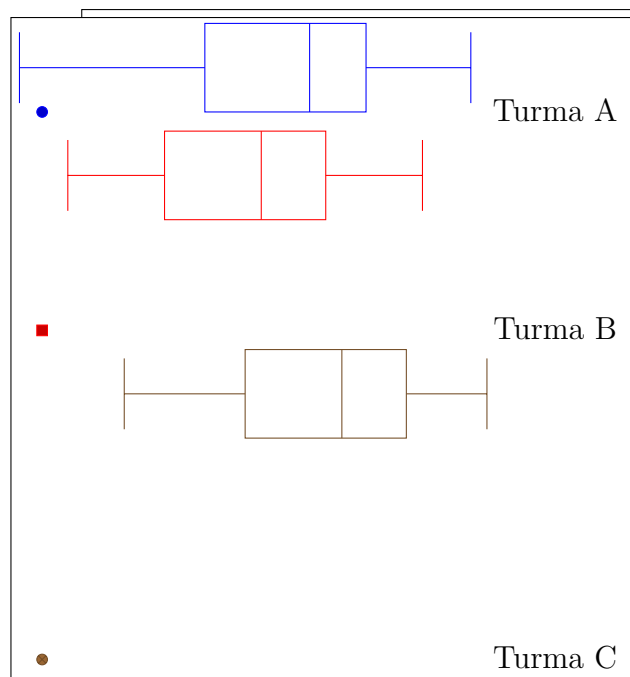


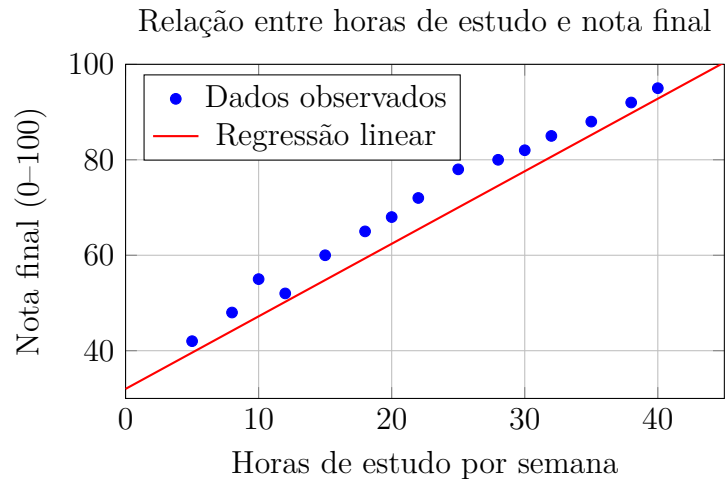
Gráfico de barras categóricas

- **Boxplot:** Mostra mediana, quartis e outliers, facilitando detecção de anomalias.  
Boxplot das notas por turma



Boxplot comparativo de três turmas

- **Gráfico de Dispersão:** Para relações entre duas variáveis quantitativas,



indicando correlações.

Gráfico de dispersão com linha de tendência

Essas ferramentas visuais resumem dados complexos, tornando-os acessíveis. Em pesquisas sociais, um histograma pode destacar desigualdades em distribuições de renda.

### 2.1.5 Assimetria e Curtose

Para uma descrição completa, avaliamos a forma da distribuição:

- **Assimetria (Skewness):** Mede o desvio da simetria. Positiva (cauda direita longa), negativa (cauda esquerda) ou zero (simétrica).

$$Skew = \frac{3(\bar{x} - Mediana)}{s} \quad (2.3)$$

- **Curtose:** Indica o achatamento das caudas. Alta curtose significa caudas pesadas (mais outliers); baixa, caudas leves.

Essas medidas ajudam a escolher testes estatísticos apropriados e detectar desvios de normalidade.

## 3 Análise Exploratória de Dados

### 3.1 Conceito básico

A Análise Exploratória de Dados (Exploratory Data Analysis, ou EDA, na sigla em inglês) representa um pilar fundamental na ciência de dados e na estatística aplicada, servindo como etapa inicial e indispensável para compreender conjuntos de dados complexos antes de proceder a modelagens mais avançadas. A EDA enfatiza a importância de investigar os dados de forma aberta e iterativa, sem pressuposições rígidas, para revelar padrões, anomalias e relacionamentos subjacentes. Diferentemente da análise confirmatória, que testa hipóteses pré-definidas, a EDA adota uma perspectiva indutiva, guiada pela curiosidade empírica e pela robustez estatística.

Deve se entender que a EDA não é apenas como uma coleção de técnicas descritivas, mas como um *framework* metodológico ancorado em princípios probabilísticos e inferenciais. Em um contexto onde os dados crescem exponencialmente (impulsionados por

avanços em *Big Data* e *Machine Learning*), a EDA assume um papel crucial na mitigação de vieses, na detecção de erros e na formulação de hipóteses informadas.

### 3.2 Fundamentos Teóricos da Análise Exploratória de Dados

A EDA baseia-se em premissas estatísticas que priorizam a resiliência dos dados a suposições paramétricas. Em via de regra, não necessariamente os dados irão aderir a distribuições ideais, como a normal, e que métodos resistentes (isto é, aqueles menos sensíveis a outliers ou violações de normalidade) devem ser privilegiados. Por exemplo, em vez de depender exclusivamente da média aritmética, que é vulnerável a valores extremos, a EDA promove o uso da mediana como medida de tendência central mais robusta.

Do ponto de vista probabilístico, a EDA incorpora conceitos de distribuições empíricas e funções de distribuição cumulativa (CDF). A CDF empírica, definida como:

$$\hat{F}_{(x)} = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x) \quad (3.1)$$

onde  $I$  é a função indicadora e  $n$  o tamanho da amostra, permite uma representação não paramétrica da distribuição dos dados, facilitando comparações com distribuições teóricas via testes como o de Kolmogorov-Smirnov. Essa abordagem estatística permite identificar desvios sistemáticos, como assimetria (skewness) ou curtose (kurtosis), quantificados respectivamente por:

$$\gamma_1 = \frac{\mu_3}{\sigma^3} \text{ e } \gamma_2 = \frac{\mu_4}{\sigma^4} - 3 \quad (3.2)$$

onde  $\mu_k$  são os momentos centrais e  $\sigma$  o desvio padrão, respectivamente.

Além disso, a EDA integra princípios de inferência bayesiana informal, onde priors subjetivos são atualizados com evidências dos dados. Embora não formalize modelos bayesianos completos, a iteração exploratória refina crenças iniciais sobre a estrutura dos dados, alinhando-se à filosofia de Tukey de "detecção de surpresas". Em resumo, os fundamentos teóricos da EDA com viés estatístico enfatizam a robustez, a não parametricidade e a iteração, preparando o terreno para análises mais rigorosas.

### 3.3 Técnicas Estatísticas Centrais na EDA

As técnicas estatísticas formam o cerne da EDA, divididas em univariadas, bivariadas e multivariadas. Na análise univariada, o foco recai sobre a distribuição individual de variáveis. Medidas de tendência central — média ( $\bar{x} = \frac{1}{n} \sum x_i$ ), mediana e moda — são complementadas por medidas de dispersão, como variância ( $s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$ ), desvio padrão e intervalo interquartil (IQR = Q3 - Q1). O IQR, em particular, é estatisticamente robusto para detectar outliers via regra de Tukey: valores abaixo de  $Q1 - 1.5 \times \text{IQR}$  ou acima de  $Q3 + 1.5 \times \text{IQR}$  são considerados anômalos.

A análise de distribuições envolve testes de normalidade, como o Shapiro-Wilk, que testa a hipótese nula de que os dados seguem uma distribuição normal:

$$W = \frac{(\sum a_i x_{(i)})^2}{\sum (x_i - \bar{x})^2} \quad (3.3)$$

onde  $a_i$  são coeficientes tabulados. Se rejeitada, distribuições alternativas (e.g., log-normal, Poisson) são exploradas. Para dados categóricos, frequências relativas e testes

qui-quadrados revelam associações inesperadas:

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i} \quad (3.4)$$

Na análise bivariada, coeficientes de correlação são pivôs. O coeficiente de Pearson:

$$cr = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}} \quad (3.5)$$

assume linearidade e normalidade, enquanto o de Spearman, baseado em ranks, é não paramétrico e robusto a outliers. Para variáveis categóricas, o coeficiente de contingência ou o teste de independência qui-quadrado é aplicado. Análises multivariadas estendem isso via matrizes de correlação e análise de componentes principais (PCA), onde autovalores e autovetores decompõem a variância:  $\mathbf{X} = \mathbf{TP}^T + \mathbf{E}$ , reduzindo dimensionalidade e identificando clusters latentes.

### 3.4 Ferramentas Visuais e Computacionais com Ênfase Estatística

A visualização é o braço direito da EDA, transformando abstrações estatísticas em insights intuitivos. Histogramas e gráficos de densidade kernel (KDE) ilustram distribuições, com o KDE calculado pela expressão abaixo:

$$\hat{f}(x) = \frac{1}{nh} \sum K\left(\frac{x - x_i}{h}\right) \quad (3.6)$$

onde  $K$  é o kernel e  $h$  o bandwidth. Boxplots, introduzidos por Tukey, resumem quartis e outliers, facilitando comparações entre grupos via testes não paramétricos.

Scatterplots e heatmaps de correlação visualizam relacionamentos bivariados, com linhas de regressão linear:

$$\hat{y} = b_0 + b_1x, \text{ onde } b_1 = r \frac{s_y}{s_x} \quad (3.7)$$

destacando tendências. Para dados multivariados, pairplots e gráficos de coordenadas paralelas revelam interações complexas. Em contextos estatísticos avançados, Q-Q plots comparam quantis empíricos com teóricos, diagnosticando desvios de normalidade.

Computacionalmente, linguagens como **R** (com pacotes como *ggplot2* e *exploratory*) e Python (com *pandas*, *seaborn* e *statsmodels*) operacionalizam essas técnicas. Por exemplo, em Python, a função **pandas.describe()** gera sumários estatísticos, enquanto **scipy.stats** executa testes. A integração com machine learning estende a EDA para descoberta de padrões não supervisionados, mantendo o viés estatístico via validação cruzada e métricas como o coeficiente de silhueta.



## 4 Amostragem

### 4.1 Conceito básico

A amostragem constitui um dos pilares fundamentais da estatística inferencial, permitindo que conclusões sobre uma população inteira sejam extraídas a partir do estudo de uma subparte dela, denominada amostra. Em contextos onde o censo — a análise exaustiva de todos os elementos da população — é impraticável devido a limitações de tempo, custo ou recursos, a amostragem emerge como uma ferramenta indispensável. Conceitualmente, a população refere-se ao conjunto completo de indivíduos ou unidades de interesse, caracterizada por parâmetros como média  $\mu$  e variância  $\sigma^2$ , enquanto a amostra é um subconjunto selecionado, cujas estatísticas (e.g., média amostral  $\bar{x}$  e variância  $s^2$ ) servem como estimadores desses parâmetros.

A teoria da amostragem enfatiza a importância de métodos que minimizem vieses e erros amostrais, garantindo a representatividade e a validade inferencial. Distinguem-se dois grandes paradigmas: a amostragem probabilística, onde cada elemento da população tem probabilidade conhecida e não nula de ser selecionado, permitindo o cálculo de erros padrão e intervalos de confiança; e a amostragem não probabilística, baseada em critérios subjetivos, mais suscetível a vieses, mas útil em pesquisas exploratórias ou com populações de difícil acesso.

### 4.2 Fundamentação Teórica da Amostragem

A amostragem probabilística baseia-se no princípio de que a seleção aleatória reduz o viés de seleção, permitindo que a distribuição amostral de uma estatística siga leis probabilísticas conhecidas, como o Teorema Central do Limite (TCL). Pelo TCL, para amostras suficientemente grandes ( $n > 30$ ), a média amostral  $\bar{x}$  aproxima-se de uma distribuição normal com média  $\mu$  e variância  $\sigma^2/n$ , independentemente da distribuição populacional.

O erro amostral, quantificado pelo erro padrão, é dado por  $SE = \sigma/\sqrt{n}$  para a média, com correção para populações finitas:

$$SE = \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}} \quad (4.1)$$

onde  $N$  é o tamanho populacional. Já o viés ocorre quando o método de seleção sistematicamente super ou subestima o parâmetro. Em amostragens não probabilísticas, o viés é incontornável, impedindo generalizações rigorosas. Outro conceito chave é o quadro amostral (sampling frame), lista completa e atualizada da população, essencial para evitar erros de cobertura.

### 4.3 Métodos de Amostragem Probabilística

Os métodos probabilísticos garantem inferência estatística válida. São eles:

- **Amostragem Aleatória Simples (AAS):** Cada elemento tem probabilidade igual de seleção ( $p = n/N$ ). Pode ser com ou sem reposição. Fórmula de variância da média:

$$Var(\bar{x}) = \frac{\sigma^2}{n} \left(1 - \frac{n}{N}\right) \quad (4.2)$$

Simples e sem viés, mas ineficiente para populações heterogêneas.

- **Amostragem Sistemática:** Seleciona-se um elemento inicial aleatório  $k$  e, em seguida, cada  $i$ -ésimo ( $i = N/n$ ). Eficiente operacionalmente, mas suscetível a periodicidades na lista.
- **Amostragem Estratificada:** A população é dividida em estratos homogêneos (e.g., por idade ou renda), com seleção aleatória proporcional ou ótima dentro de cada. Reduz a expressão da variância por:

$$Var(\bar{x}_{est}) = \sum W_h^2 \frac{\sigma_h^2}{n_h} \quad (4.3)$$

onde  $W_h$  é o peso do estrato. Ideal para heterogeneidade conhecida.

- **Amostragem por Conglomerados (Cluster):** Divide-se em clusters heterogêneos (e.g., escolas), selecionando clusters aleatoriamente e amostrando todos ou subamostra dentro. Útil para populações dispersas, mas aumenta variância se clusters forem semelhantes internamente.

#### 4.4 Métodos de Amostragem Não Probabilística

Embora menos rigorosos, são comuns em pesquisas qualitativas:

- **Amostragem por Conveniência:** Seleção de elementos acessíveis. Rápida, mas propensa a viés.
- **Amostragem por Cotas:** Similar à estratificada, mas sem aleatoriedade; fixa cotas por características.
- **Amostragem Intencional ou por Julgamento:** É onde o especialista seleciona elementos típicos.
- **Amostragem "Bola de Neve":** Indicada para populações raras; participantes indicam outros.

Esses métodos não permitem cálculo de erros padrão, limitando a generalização.

#### 4.5 Dimensionamento Amostral

O tamanho da amostra  $n$  é determinado por fórmulas que equilibram precisão e custo. Para estimar uma média com margem de erro  $e$  e nível de confiança  $1 - \alpha$  (tipicamente 95%,  $z = 1.96$ ):

$$n = \frac{z^2 \sigma^2}{e^2} \quad (4.4)$$

Para populações finitas:

$$n = \frac{n_0}{1 + \frac{n_0 - 1}{N}} \quad (4.5)$$

onde  $n_0$  é o  $n$  infinito.

Para proporções ( $p$  estimada, frequentemente 0.5 para máxima variância):

$$n = \frac{z^2 p(1 - p)}{e^2} \quad (4.6)$$

Fatores como poder estatístico (em testes de hipóteses) e efeito esperado também influenciam.

## 5 Probabilidade

### 5.1 Conceito Básico

A teoria da probabilidade constitui o alicerce matemático para o estudo de fenômenos aleatórios, fornecendo ferramentas rigorosas para quantificar incerteza e modelar eventos incertos. Surgida no século XVII com contribuições de Blaise Pascal e Pierre de Fermat no contexto de jogos de azar, e formalizada por Andrey Kolmogorov em 1933 por meio de uma abordagem axiomática, a probabilidade transcende sua origem lúdica para se tornar indispensável em campos como estatística inferencial, física quântica, inteligência artificial e economia.

Do ponto de vista matemático, a probabilidade é uma medida sobre espaços de eventos, satisfazendo axiomas que garantem consistência lógica. Em estatística, ela sustenta a inferência, permitindo que amostras informem sobre populações via distribuições probabilísticas. Este texto explora os fundamentos axiomáticos, conceitos chave como probabilidade condicional e independência, distribuições de probabilidade, teoremas limitantes centrais e aplicações, enfatizando o rigor matemático e as implicações para a análise de dados.

### 5.2 Fundamentos e Definições da Probabilidade

Na literatura podemos encontrar três definições de probabilidade: **Clássica**, **Geométrica** e **Frequentista**.

Na definição clássica todos os elementos de um espaço amostral possuem a mesma chance de acontecerem. Seja um evento  $A$  de interesse, associado a um espaço amostral  $\Omega$ . Então a probabilidade de ocorrência do evento  $A$ , será a razão entre o número de elementos do evento de interesse com o número de elementos do espaço amostral.

$$P(A) = \frac{n(A)}{n(\Omega)} \quad (5.1)$$

Considerando que o espaço amostral pode ser *não enumerável*, então o conceito de probabilidades se aplicará ao comprimento de intervalos, medida de áreas ou similares, dando origem a probabilidade geométrica.

$$P(A) = \frac{\text{comprimento de } A}{\text{comprimento de } \Omega} \quad (5.2)$$

Por fim, na definição frequentista, deve-se considerar o limite das frequências relativas como o valor da probabilidade. Assim, seja  $n_A$  o número de ocorrências de  $A$  em  $n$  repetições independentes do experimento em questão, temos:

$$P(A) = \lim_{n \rightarrow \infty} \frac{n_A}{n} \quad (5.3)$$

Em síntese, isso significa que a medida que são realizados os experimentos, a probabilidade de ocorrência de um evento determinado se aproxima do verdadeiro valor a medida que o número de realizações tende ao infinito.

Para isso, considere um experimento que consiste em lançar uma moeda 10, 50, 100 e 1000 vezes, e observar o número de ensaios em que o resultado é cara. Os resultados podem ser verificados na Figura 01 abaixo. Observe que a medida que o número de ensaios cresce

a probabilidade acumulada da ocorrência de cara converge para sua verdadeira ocorrência, isto é,  $P(A) = 0,5$ .

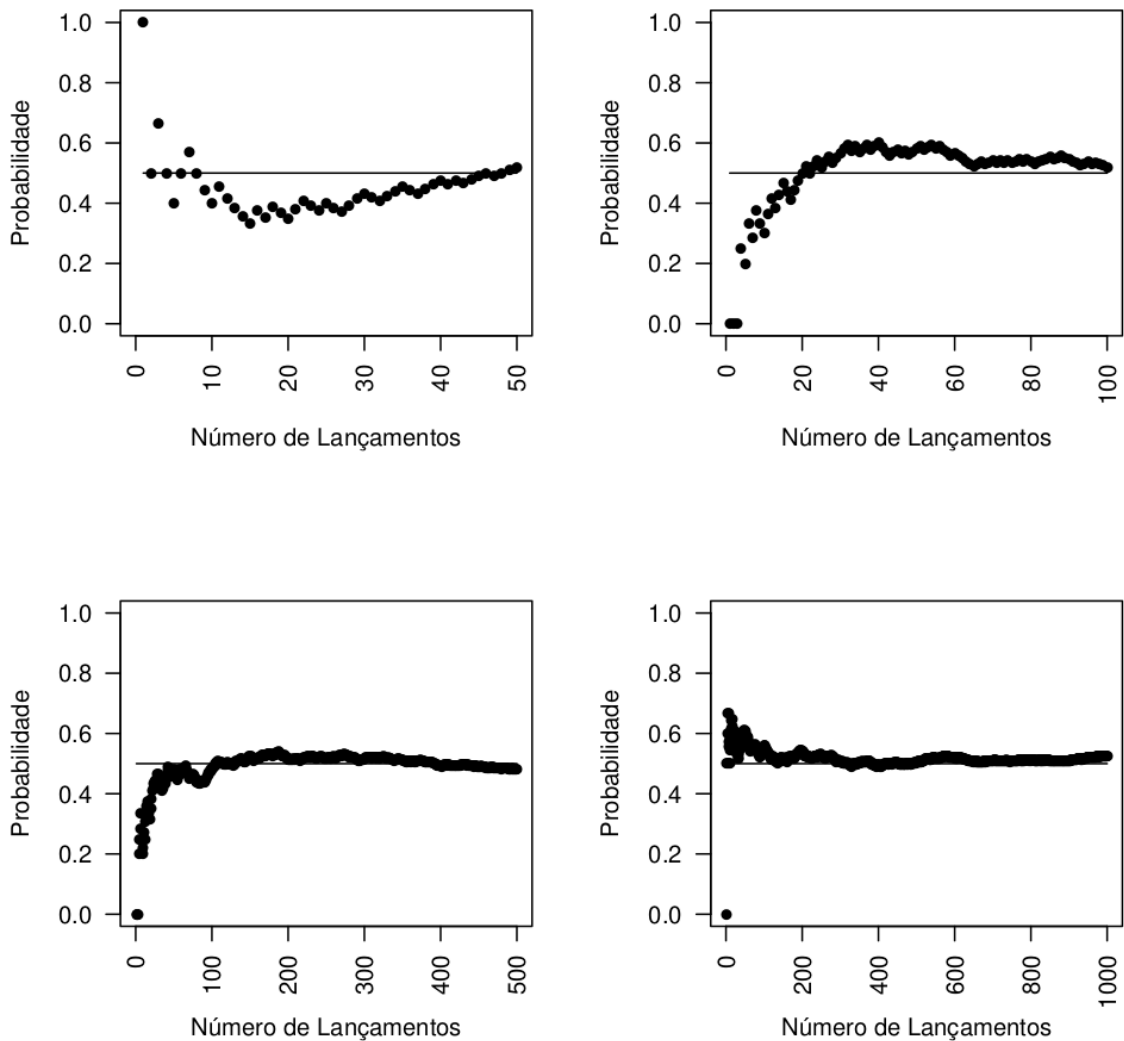


Figura 1: Convergência da frequência relativa de caras ao lançar uma moeda várias vezes.

Considerando que as definições anteriores são úteis para o cálculo de diversos problemas práticos e teóricos, é necessário enunciar uma série de axiomas para que se tenha uma formulação mais rigorosa para o conceito de probabilidade. Por volta do ano de 1930 o russo A. N. Kolmogorov apresentou esses axiomas matemáticos para definir probabilidade.

- **Axioma 1:** Para qualquer evento  $A$ ,  $P(A) \geq 0$ .

O **Axioma 1** de Kolmogorov é assumido como verdadeiro por definição. Não há prova derivada; é um postulado fundamental que garante que probabilidades não sejam negativas.

- **Axioma 2:** Seja  $\Omega$  o espaço amostral associado ao experimento aleatório, então  $P(\Omega) = 1$ .

O **Axioma 2** de Kolmogorov também é assumido por definição, pois representa a certeza de que algo no espaço amostral ocorre.

- **Axioma 3:** Se  $A_1, A_2, \dots, A_k$  for um conjunto finito de eventos mutuamente exclusivos, isto é,  $A_i \cap A_j = \emptyset$ , então  $P(\bigcup_{i=1}^k A_i) = \sum_{i=1}^k P(A_i)$ .

O **Axioma 3** de Kolmogorov também é assumido por definição, pois para coleções finitas, reduz-se a  $P(E_1 \cup \dots \cup E_n) = P(E_1) + \dots + P(E_n)$  quando disjuntos.

Ou seja, No experimento aleatório de lançar uma moeda e observar a face superior, o espaço amostral é  $\Omega = H, T$ . Pelos axiomas citados anteriormente, tem-se que  $P(\Omega) = 1$ . Logo:

$$1 = P(\Omega) = P(H \cup T) = P(H) + P(T) = 0,5 + 0,5 \quad (5.4)$$

### 5.3 Propriedades Derivadas

As propriedades abaixo são teoremas provados a partir dos axiomas supracitados.

#### 1. Probabilidade do conjunto vazio:

$$P(\emptyset) = 0.$$

*Demonstração:* Note que  $\emptyset \cup \emptyset = \emptyset$ . Pelo Axioma 3,

$$P(\emptyset) = P(\emptyset \cup \emptyset) = P(\emptyset) + P(\emptyset) \implies P(\emptyset) = 0.$$

#### 2. Probabilidade do complemento: Para qualquer evento $A$ ,

$$P(A^c) = 1 - P(A).$$

*Demonstração:* Temos  $A \cup A^c = \Omega$  e  $A \cap A^c = \emptyset$ . Pelo Axioma 3 e Axioma 2,

$$P(\Omega) = P(A \cup A^c) = P(A) + P(A^c) = 1 \implies P(A^c) = 1 - P(A).$$

#### 3. Monotonicidade: Se $A \subseteq B$ , então

$$P(A) \leq P(B).$$

*Demonstração:* Se  $A \subseteq B$ , então  $B = A \cup (B \setminus A)$  com  $A \cap (B \setminus A) = \emptyset$ . Pelo Axioma 3,

$$P(B) = P(A) + P(B \setminus A).$$

Como  $P(B \setminus A) \geq 0$  (Axioma 1), segue que  $P(B) \geq P(A)$ .

#### 4. Limites da probabilidade: Para qualquer evento $A$ ,

$$0 \leq P(A) \leq 1.$$

*Demonstração:*

- $P(A) \geq 0$  vem do Axioma 1.
- Como  $A \subseteq \Omega$ , pela monotonicidade  $P(A) \leq P(\Omega) = 1$ . Ou ainda:  $P(A) + P(A^c) = 1$  e  $P(A^c) \geq 0 \implies P(A) \leq 1$ .

**5. Princípio de inclusão-exclusão para dois eventos:**

$$P(A \cup B) = P(A) + P(B) - P(A \cap B).$$

*Demonstração:* Escreva  $A \cup B = A \cup (B \setminus A)$ , onde  $B \setminus A = B \cap A^c$  e os conjuntos são disjuntos. Pelo Axioma 3,

$$P(A \cup B) = P(A) + P(B \setminus A).$$

Agora,  $B = (A \cap B) \cup (B \setminus A)$  com interseção vazia, logo

$$P(B) = P(A \cap B) + P(B \setminus A) \implies P(B \setminus A) = P(B) - P(A \cap B).$$

Substituindo,

$$P(A \cup B) = P(A) + P(B) - P(A \cap B).$$

**6. Probabilidade da união de eventos disjuntos (caso finito):** Se  $A \cap B = \emptyset$ ,

$$P(A \cup B) = P(A) + P(B).$$

*Demonstração:* Caso particular do Axioma 3 ou do princípio de inclusão-exclusão:

$$P(A \cup B) = P(A) + P(B) - P(\emptyset) = P(A) + P(B).$$

Essas são as propriedades mais fundamentais e comumente usadas na probabilidade clássica. Elas valem tanto para espaços de probabilidade finitos quanto contínuos, desde que os axiomas de Kolmogorov sejam satisfeitos.

## 6 Distribuição Normal

A distribuição normal, frequentemente chamada de curva de sino ou distribuição gaussiana, é uma das ferramentas mais fundamentais e versáteis da estatística e da probabilidade. Ela descreve como os valores de uma variável se distribuem ao redor de um ponto central, com a maioria dos dados concentrados no meio e diminuindo simetricamente para os lados. Imagine uma montanha simétrica, onde o pico representa o valor mais comum, e as encostas mostram como os valores se tornam menos frequentes à medida que se afastam do centro. Essa forma é tão comum na natureza e nos dados humanos que a distribuição normal se tornou um pilar para entender padrões em tudo, desde alturas de pessoas até erros de medição em experimentos científicos.

### 6.1 Origem

A distribuição normal tem raízes no século XVIII, quando matemáticos como Abraham de Moivre começaram a estudá-la como uma aproximação para distribuições binomiais em jogos de azar. Por exemplo, ao lançar uma moeda muitas vezes, a probabilidade de obter um número de caras próximo à metade do total se concentra em torno de um valor médio, formando algo semelhante a uma curva de sino. No início do século XIX, Pierre-Simon Laplace e Carl Friedrich Gauss refinaram essas ideias, aplicando-as a erros de observação em astronomia. Gauss, em particular, usou a distribuição para modelar variações em medições celestes, o que levou à sua associação com o nome "gaussiana".

Ao longo do século XIX, figuras como Adolphe Quetelet aplicaram o conceito a dados sociais, como alturas e pesos de populações humanas, demonstrando que muitos traços biológicos seguem padrões normais. Francis Galton, um pioneiro na estatística, admirava tanto essa distribuição que a descreveu como uma "lei de frequência de erro" que revelava uma ordem cósmica. No século XX, com o avanço da estatística inferencial por Ronald Fisher e outros, a normal se consolidou como base para testes estatísticos e modelagens. Hoje, em uma era de big data e inteligência artificial, ela continua relevante, adaptada a contextos digitais e computacionais.

Essa evolução reflete uma transição de curiosidades matemáticas para uma ferramenta essencial na ciência empírica. A distribuição normal não foi "inventada", mas descoberta como um padrão recorrente, destacando como a matemática pode capturar a essência da variabilidade no mundo.

## 6.2 Definição Formal e Função Densidade

Seja  $X$  uma variável aleatória contínua. Dizemos que  $X$  segue uma **distribuição normal** com média  $\mu$  e variância  $\sigma^2 > 0$  (denotada  $X \sim \mathcal{N}(\mu, \sigma^2)$ ) quando sua função densidade de probabilidade é:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right], \quad -\infty < x < +\infty \quad (6.1)$$

onde:

- $\mu$  (média)  $\rightarrow$  centro da curva (pico do sino)
- $\sigma$  (desvio padrão)  $\rightarrow$  controla a largura do sino
- $\sigma^2$  (variância)  $\rightarrow$  medida de dispersão ao quadrado

## 6.3 Distribuição Normal Padrão

A normal padrão é o caso particular em que  $\mu = 0$  e  $\sigma = 1$ . Denotamos  $Z \sim \mathcal{N}(0, 1)$ . Sua densidade é:

$$\phi(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) \quad (6.2)$$

Qualquer normal pode ser padronizada pela transformação:

$$Z = \frac{X - \mu}{\sigma} \quad (6.3)$$

## 6.4 Principais Propriedades

- Simetria perfeita:  $f(\mu + t) = f(\mu - t)$  para todo  $t$
- Média = Mediana = Moda =  $\mu$
- Variância =  $\sigma^2$  e desvio padrão =  $\sigma$
- Assimetria (skewness) = 0

- Curtose excessiva = 0 (curva mesocúrtica)
- Função geradora de momentos:

$$M(t) = \exp\left(\mu t + \frac{\sigma^2 t^2}{2}\right)$$

- **Reprodutividade (fechamento sob combinações lineares):**  
Se  $X_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$  forem independentes e  $a_i, b$  constantes, então

$$\sum_i a_i X_i + b \sim \mathcal{N}\left(\sum_i a_i \mu_i + b, \sum_i a_i^2 \sigma_i^2\right)$$

## 6.5 Regra Empírica (68–95–99,7)

$$\begin{aligned} P(|X - \mu| < \sigma) &\approx 0,683 \\ P(|X - \mu| < 2\sigma) &\approx 0,954 \\ P(|X - \mu| < 3\sigma) &\approx 0,997 \end{aligned}$$

Na normal padrão:

$$\begin{aligned} P(-1 < Z < 1) &\approx 68,3\% \\ P(-2 < Z < 2) &\approx 95,4\% \\ P(-3 < Z < 3) &\approx 99,7\% \end{aligned}$$

## 6.6 Função de Distribuição Cumulativa

$$\Phi(z) = P(Z \leq z) = \int_{-\infty}^z \phi(t) dt$$

Não possui forma fechada, mas está extensivamente tabelada e implementada em softwares.

## 7 Teorema Central do Limite (TCL)

O Teorema Central do Limite (TCL) é um dos resultados mais profundos e influentes da teoria da probabilidade e da estatística. Ele afirma que, sob condições relativamente brandas, a distribuição da soma (ou da média) de um grande número de variáveis aleatórias independentes e identicamente distribuídas (i.i.d.) tende a uma distribuição normal, independentemente da forma da distribuição original das variáveis individuais. Essa convergência ocorre à medida que o tamanho da amostra aumenta, explicando por que a distribuição normal é tão ubiquamente observada em fenômenos reais.

Formalizado no século XX, embora com raízes no trabalho de Abraham de Moivre (1733) e Pierre-Simon Laplace (1810), o TCL moderno deve muito a contribuições de matemáticos como Aleksandr Lyapunov e Jarl Waldemar Lindeberg. Sua importância reside no fato de que ele fornece a base teórica para grande parte da inferência estatística clássica, permitindo o uso de métodos paramétricos (como testes de hipóteses e intervalos



de confiança) mesmo quando a distribuição populacional subjacente é desconhecida ou não normal.

Outros resultados semelhantes, como aproximações para distribuições binomiais em grandes amostras, reforçam essa centralidade. Em resumo, a normal não é imposta; ela emerge organicamente da soma de influências aleatórias, refletindo a complexidade do mundo real.

## 7.1 Definição

Considere  $X_1, X_2, \dots, X_n$  variáveis aleatórias independentes e identicamente distribuídas com média  $\mathbb{E}[X_i] = \mu$  e variância  $\text{Var}(X_i) = \sigma^2 < \infty$ . Defina a média amostral  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ . O TCL clássico (versão de Lindeberg-Lévy) afirma que, quando  $n \rightarrow \infty$ :

$$\sqrt{n} (\bar{X}_n - \mu) / \sigma \xrightarrow{d} \mathcal{N}(0, 1) \quad (7.1)$$

onde  $\xrightarrow{d}$  denota convergência em distribuição e  $\mathcal{N}(0, 1)$  é a distribuição normal padrão. Equivalentemente, para a soma  $S_n = \sum_{i=1}^n X_i$ :

$$\frac{S_n - n\mu}{\sigma\sqrt{n}} \xrightarrow{d} \mathcal{N}(0, 1) \quad (7.2)$$

## 7.2 Condições de Validade e Versões Mais Gerais

A versão mais simples requer que as variáveis sejam i.i.d. com variância finita. Versões mais gerais relaxam essas hipóteses:

- **Independência, mas não idêntica distribuição** (Lindeberg): Desde que nenhuma variável domine a variância total.
- **Dependência fraca**: para séries temporais ou processos estacionários
- **Momentos finitos** algumas versões exigem apenas a existência de momentos de ordem  $2 + \varepsilon$

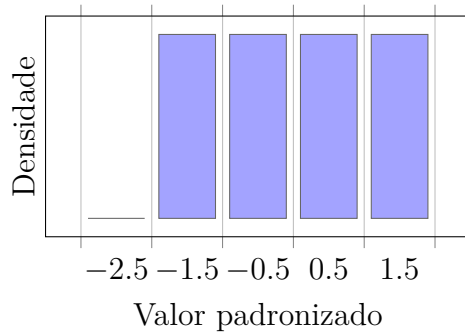
O TCL falha se a variância for infinita (ex.: distribuição de Cauchy), onde a média amostral não converge para uma normal.

Na prática, para amostras com  $n \geq 30$ , a aproximação já é razoável na maioria dos casos, especialmente se a distribuição original não for muito assimétrica.

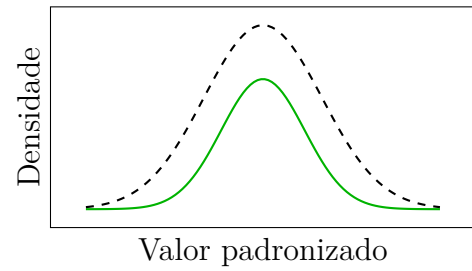
## 7.3 Interpretação Intuitiva

Intuitivamente, o TCL reflete o fato de que a soma de muitos efeitos aleatórios independentes tende a "equilibrar-se". Erros positivos e negativos se cancelam, e a variabilidade relativa diminui com  $\sqrt{n}$ . Por exemplo, ao lançar um dado muitas vezes, a média tende a 3,5, com distribuição aproximadamente normal, mesmo que a distribuição uniforme original seja discreta e plana.

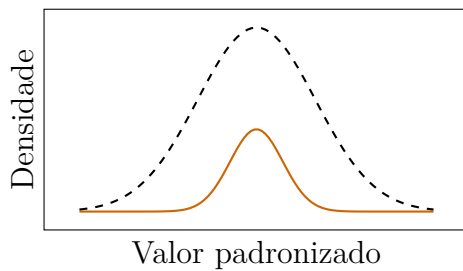
$n = 1$  (distribuição original - exemplo uniforme)



$n = 2$



$n = 5$



$n = 25$

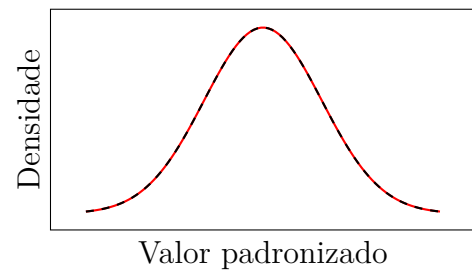


Figura 2: Ilustração do Teorema Central do Limite: convergência da distribuição da média amostral padronizada para a distribuição normal padrão  $\mathcal{N}(0, 1)$  à medida que o tamanho da amostra  $n$  aumenta. A curva tracejada preta representa a normal limite em todos os painéis.

O TCL é o fundamento para:

- Intervalos de confiança para médias:  $\bar{X} \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$  (ou com estimativa  $s$  para  $t$  de Student em  $n$  pequeno).
- Testes de hipóteses paramétricos (z-test, t-test).
- Aproximações normais para outras estatísticas (ex.: proporções com  $n$  grande).
- Modelos lineares e regressão, onde erros assumem normalidade aproximada.

Sem o TCL, a inferência estatística seria limitada a distribuições conhecidas ou métodos não paramétricos mais conservadores.

## 8 Aplicações na Análise e Ciência de Dados

A análise de dados e a ciência de dados (data science) representam campos interdisciplinares que integram conceitos estatísticos fundamentais para extrair insights acionáveis de conjuntos de dados complexos. Tópicos como Análise Exploratória de Dados (EDA), amostragem, probabilidade, distribuição normal e o Teorema Central do Limite (TCL) formam o cerne metodológico dessas disciplinas. Eles permitem não apenas a compreensão inicial dos dados, mas também a construção de modelos preditivos robustos, a mitigação de incertezas e a tomada de decisões informadas em contextos reais, como negócios, saúde, finanças e inteligência artificial.

Em um mundo impulsionado por big data, esses conceitos estatísticos são aplicados para lidar com volumes massivos de informações, identificar padrões ocultos e validar hipóteses. Por exemplo, em machine learning, eles sustentam algoritmos que vão desde regressões lineares até redes neurais profundas. Este texto explora as aplicações práticas desses tópicos, ancoradas em exemplos contemporâneos, destacando sua relevância em cenários de análise de dados e ciência de dados.

## 8.1 Aplicações da Análise Exploratória de Dados (EDA)

A Análise Exploratória de Dados é a etapa inicial em qualquer projeto de data science, envolvendo técnicas descritivas, visuais e estatísticas para compreender a estrutura, padrões e anomalias nos dados. Sua aplicação vai além da mera descrição: ela guia a limpeza de dados, a feature engineering e a formulação de hipóteses para modelagens subsequentes.

Em contextos empresariais, a EDA é crucial para segmentação de clientes e detecção de anomalias. Por exemplo, em análise de mercado, técnicas como gráficos de dispersão e histogramas ajudam a identificar tendências em dados de vendas, permitindo otimizar estratégias de precificação. Na saúde, EDA é usada para explorar conjuntos de dados epidemiológicos, revelando correlações entre variáveis como idade e incidência de doenças, o que auxilia em modelagens preditivas para surtos. Visualizações são centrais na EDA, como boxplots para detecção de outliers ou heatmaps para correlações.

Essas ferramentas, implementadas em bibliotecas como **pandas** e **seaborn** em Python, facilitam a identificação de vieses em datasets, essencial para ética em IA. Em projetos de data science, uma EDA robusta pode reduzir o tempo de desenvolvimento de modelos em até 30%, ao evitar suposições errôneas.

## 8.2 Aplicações da Técnica de Amostragem

A amostragem é essencial em data science para lidar com datasets massivos, onde processar o todo é computacionalmente inviável. Técnicas probabilísticas, como amostragem aleatória simples ou estratificada, garantem representatividade, reduzindo custos e tempo de análise enquanto minimizam vieses.

Em big data analytics, amostragem é aplicada em processamentos distribuídos, como no Hadoop ou Spark, para estimar estatísticas populacionais a partir de subconjuntos. Por exemplo, em pesquisas de opinião ou análise de redes sociais, amostragem por conglomerados é usada para inferir comportamentos de usuários a partir de amostras de posts ou perfis. Na ciência ambiental, técnicas como amostragem sistemática auxiliam na modelagem de dados climáticos, permitindo previsões globais baseadas em dados locais.

## 8.3 Aplicações da Teoria da Probabilidade

A probabilidade fornece o framework para lidar com incerteza em data science, permitindo quantificar riscos e fazer inferências bayesianas. Conceitos como distribuições probabilísticas e independência são aplicados em modelagens estocásticas e algoritmos de aprendizado.

Em análise preditiva, probabilidade é usada em regressão logística para estimar probabilidades de eventos binários, como churn de clientes. Na IA, redes bayesianas modelam dependências em dados, aplicadas em diagnósticos médicos ou recomendadores de conteúdo.