

# Aula 2 – Revisão de Conceitos Essenciais de Estatística

Bem-vindo(a) à segunda etapa da sua jornada no Curso de Modelagem Preditiva Avançada! Antes de mergulharmos nas complexidades dos algoritmos e das arquiteturas de modelos, é fundamental solidificarmos as bases. Pense na estatística como o alicerce de um arranha-céu: sem uma fundação robusta, a estrutura, por mais imponente que seja, corre o risco de desabar. Da mesma forma, sem uma compreensão clara dos conceitos estatísticos, a construção de modelos preditivos eficazes e, mais importante, interpretáveis, torna-se um desafio.

Nesta aula, não vamos apenas revisar fórmulas, mas sim desvendar a lógica por trás delas, conectando cada conceito à sua aplicação prática no universo da modelagem preditiva e da análise de dados. Nosso objetivo é que você não apenas lembre o que é uma média ou um desvio padrão, mas que compreenda *por que* eles são importantes, *quando* usá-los e *como* eles influenciam as decisões que você tomará ao construir seus modelos. Ao final desta sessão, você estará apto(a) a identificar e aplicar as principais medidas estatísticas, entender a natureza das distribuições de dados e interpretar os resultados de testes de hipóteses, elementos cruciais para qualquer cientista de dados ou analista que busca ir além do "apertar botões".

Prepare-se para uma revisão que transformará conceitos abstratos em ferramentas poderosas. Vamos explorar as medidas de tendência central e dispersão, mergulhar nas distribuições de probabilidade mais comuns e desmistificar os testes de hipóteses e o p-valor, tudo isso com um olhar voltado para as tendências atuais como a Inteligência Artificial Explicável (XAI).

# A Essência dos Dados: Onde Estamos e Para Onde Vamos?

📄 **Ponto de Partida:** Transformar dados brutos em insights acionáveis começa com medidas de tendência central.

Imagine que você está diante de uma montanha de dados brutos, talvez registros de vendas, informações de clientes ou leituras de sensores. Essa massa de números, por si só, não diz muita coisa. É como ter todas as peças de um quebra-cabeça espalhadas na mesa: você sabe que há uma imagem ali, mas não consegue vê-la. Para começar a dar sentido a essa informação, precisamos de ferramentas que nos ajudem a resumir e a descrever o que está acontecendo.

É aqui que entram as medidas de tendência central. Elas são como os pontos de referência em um mapa, indicando o "centro" ou o valor mais típico de um conjunto de dados. Entender onde a maioria dos seus dados se concentra é o primeiro passo para qualquer análise, seja para prever o comportamento do consumidor ou para otimizar um processo industrial. Sem essa compreensão básica, qualquer modelo preditivo que você construir estará operando no escuro, sem um ponto de partida claro.

Essas medidas nos permitem ter uma visão panorâmica rápida, um "primeiro olhar" sobre o comportamento geral de uma variável. Elas são a base para perguntas mais complexas, como "qual é o valor médio de um cliente?" ou "qual o tempo de resposta mais comum do nosso sistema?". Ao dominar esses conceitos, você estará apto(a) a extrair insights valiosos mesmo antes de aplicar qualquer algoritmo sofisticado de Machine Learning.

# Média, Mediana e Moda: O Coração dos Seus Dados

## Média

A soma de todos os valores dividida pelo número total de valores. É o ponto de equilíbrio de uma gangorra.

**Quando usar:** Dados numéricos sem outliers extremos

## Mediana

O valor central quando os dados estão ordenados. A pessoa no meio da fila.

**Quando usar:** Dados com outliers ou distribuições assimétricas

## Moda

O valor que aparece com maior frequência. A cor de carro mais comum no estacionamento.

**Quando usar:** Dados categóricos ou identificar o item mais popular

Quando falamos em "centro" de um conjunto de dados, geralmente pensamos na **média**. Ela é a soma de todos os valores dividida pelo número total de valores. Pense na média como o ponto de equilíbrio de uma gangorra: se você colocar todos os pesos (seus dados) em seus respectivos lugares, a média é onde você precisaria apoiar o fulcro para que a gangorra ficasse perfeitamente nivelada. É uma medida intuitiva e amplamente utilizada, mas pode ser facilmente influenciada por valores extremos, os chamados *outliers*.


A **mediana**, por outro lado, é o valor central de um conjunto de dados quando eles estão organizados em ordem crescente ou decrescente. Se você tem uma fila de pessoas classificadas por altura, a mediana é a altura da pessoa que está exatamente no meio da fila. Ela é particularmente útil quando seus dados contêm valores muito altos ou muito baixos que poderiam distorcer a média, oferecendo uma representação mais robusta do "típico". Em cenários de renda, por exemplo, a mediana é frequentemente preferida à média para evitar que salários de milionários distorçam a percepção da renda da maioria.

Por fim, a **moda** é o valor que aparece com maior frequência em um conjunto de dados. Imagine que você está contando as cores de carros em um estacionamento; a cor que aparece mais vezes é a moda. Ela é a única medida de tendência central que pode ser usada para dados categóricos (como cores, tipos de produtos) e é útil para identificar o item mais popular ou a categoria mais comum. Um e-commerce pode usar a moda para identificar o produto mais vendido em um determinado período, direcionando estratégias de estoque e marketing.

# Comparando as Medidas de Tendência Central

Essas três medidas, embora busquem o "centro", o fazem de maneiras distintas, cada uma com suas forças e fraquezas. A escolha de qual usar depende da natureza dos seus dados e do objetivo da sua análise. Por exemplo, se você está analisando o tempo de carregamento de uma página web e há alguns carregamentos extremamente lentos (outliers), a mediana pode ser uma medida mais representativa do tempo "típico" de carregamento para a maioria dos usuários do que a média.

| Conceito       | Âmbito/Aplicação                      | Base/Origem                                      | Exemplo   |
|----------------|---------------------------------------|--|---|
| <b>Média</b>   | Dados numéricos sem outliers extremos | Soma de todos os valores / Quantidade de valores | Salário médio de uma equipe (se não houver grandes discrepâncias) |
| <b>Mediana</b> | Dados numéricos com ou sem outliers   | Valor central em dados ordenados                 | Renda mediana de uma população                                    |
| <b>Moda</b>    | Dados numéricos ou categóricos        | Valor mais frequente                             | Cor de carro mais vendida   |

 **Insight Prático:** A escolha da medida de tendência central correta é um passo crucial na fase de Análise Exploratória de Dados (EDA), que precede a construção de qualquer modelo de Machine Learning. Uma compreensão equivocada do centro dos seus dados pode levar a conclusões erradas e, conseqüentemente, a modelos preditivos falhos.

# Além da Média: Entendendo a Dispersão

"Conhecer o centro não é suficiente. Precisamos entender como os dados se espalham."

Conhecer o "centro" dos seus dados é um excelente começo, mas não é o suficiente. Imagine que você está comparando dois times de basquete. Ambos têm uma altura média de 1,90m. Isso significa que eles são idênticos? Não necessariamente. Um time pode ter todos os jogadores com cerca de 1,90m, enquanto o outro pode ter jogadores muito baixos e muito altos, cuja média ainda é 1,90m. A média, sozinha, não nos diz quão "espalhados" ou "concentrados" os dados estão em torno desse centro.

É aqui que as medidas de dispersão se tornam indispensáveis. Elas nos informam sobre a variabilidade, a amplitude ou a consistência dos dados. Entender a dispersão é crucial em diversas áreas. Em finanças, por exemplo, um investimento com alto retorno médio, mas também alta dispersão, indica um risco maior. Em controle de qualidade, uma baixa dispersão em um processo produtivo significa maior consistência e menos defeitos.

01

---

## Identificar a Variabilidade

Medir o quão espalhados estão os dados

03

---

## Detectar Outliers

Identificar valores anômalos ou extremos

02

---

## Avaliar a Qualidade

Determinar a consistência e previsibilidade

04

---

## Preparar para Modelagem

Decidir sobre pré-processamento necessário

Para um cientista de dados, a dispersão é vital para entender a qualidade e a previsibilidade dos dados. Se os dados de uma variável têm alta dispersão, significa que eles variam muito, o que pode dificultar a modelagem. Por outro lado, se a dispersão é muito baixa, a variável pode não ter poder preditivo suficiente, pois quase todos os valores são iguais. As medidas de dispersão nos dão uma visão mais completa da distribuição dos dados, complementando as medidas de tendência central.

# Variância e Desvio Padrão: A Medida da Incerteza

## Variância

A **variância** é uma das medidas de dispersão mais fundamentais. Ela quantifica o quão longe, em média, cada ponto de dado está da média. Para calculá-la, pegamos a diferença de cada valor em relação à média, elevamos essa diferença ao quadrado (para eliminar valores negativos e dar mais peso a desvios maiores), somamos todos esses quadrados e dividimos pelo número de observações (ou pelo número de observações menos um, dependendo se é uma população ou amostra). O resultado é um número que nos diz o "espalhamento" médio dos dados.

**Limitação:** Sua unidade de medida é o quadrado da unidade original dos dados, dificultando a interpretação direta.

## Desvio Padrão

No entanto, a variância tem uma desvantagem: sua unidade de medida é o quadrado da unidade original dos dados, o que dificulta a interpretação direta. Por exemplo, se seus dados estão em metros, a variância estará em metros quadrados. Para resolver isso, usamos o **desvio padrão**. O desvio padrão é simplesmente a raiz quadrada da variância. Ao tirar a raiz quadrada, voltamos à unidade de medida original dos dados, tornando-o muito mais intuitivo e fácil de interpretar.

**Vantagem:** Mesma unidade dos dados originais, facilitando a interpretação.

📌 **Pense assim:** O desvio padrão é a "distância média" que os pontos de dados estão da média. Um desvio padrão pequeno indica que os dados estão agrupados perto da média, sugerindo consistência. Um desvio padrão grande, por outro lado, significa que os dados estão mais espalhados, indicando maior variabilidade.

Em Machine Learning, entender o desvio padrão de suas *features* (variáveis de entrada) é crucial para técnicas como a padronização de dados, onde transformamos os dados para que tenham média zero e desvio padrão um, o que pode melhorar o desempenho de muitos algoritmos.

# Aplicações Práticas da Variância e Desvio Padrão

A aplicação dessas medidas vai além da simples descrição. Em um contexto de Machine Learning, ao analisar a distribuição de uma variável como a idade dos clientes, um desvio padrão alto pode indicar uma base de clientes muito diversa, o que pode exigir modelos mais complexos ou a segmentação dos dados. Já um desvio padrão baixo pode sugerir que a idade não é uma variável que varia muito entre seus clientes, e talvez outras variáveis sejam mais preditivas.

## Detecção de Outliers

Valores que estão a muitos desvios padrão de distância da média são fortes candidatos a serem anomalias, que podem ser erros de coleta de dados ou eventos raros de interesse.

## Padronização de Dados

Transformar dados para média zero e desvio padrão um melhora o desempenho de algoritmos sensíveis à escala das variáveis.

## Avaliação de Consistência

Processos com baixo desvio padrão indicam maior previsibilidade e controle de qualidade.

| Conceito             | Âmbito/Aplicação                                       | Base/Origem   | Exemplo   |
|----------------------|--|---|---|
| <b>Variância</b>     | Quantifica o espalhamento dos dados em relação à média | Média dos quadrados das diferenças entre cada valor e a média | Variância da altura dos alunos em uma turma       |
| <b>Desvio Padrão</b> | Medida de dispersão na mesma unidade dos dados         | Raiz quadrada da variância                                    | Desvio padrão da temperatura diária de uma cidade |

Entender a variância e o desvio padrão permite que você avalie a confiabilidade de suas medidas de tendência central e tome decisões mais informadas sobre o pré-processamento de dados, um passo crítico antes de alimentar qualquer algoritmo de Machine Learning.

# O Mundo da Probabilidade: Previsões e Incertezas

## A probabilidade é a linguagem da incerteza

Até agora, falamos sobre descrever o que já aconteceu com nossos dados. Mas e se quisermos prever o que *pode* acontecer? É aqui que entramos no fascinante mundo da probabilidade. A probabilidade é a linguagem da incerteza, a ferramenta matemática que nos permite quantificar a chance de um evento ocorrer. Seja para prever se um cliente irá cancelar um serviço, se uma transação é fraudulenta ou se um email é spam, a probabilidade está no cerne de todas as decisões baseadas em dados e, conseqüentemente, da modelagem preditiva.

### Modelos de Classificação

Produzem probabilidades (ex: 95% de chance de ser spam)

### Modelos de Regressão

Construídos sob premissas de distribuições de probabilidade dos erros

### Quantificação de Confiança

Permite avaliar a certeza das previsões do modelo

Sem uma compreensão sólida dos conceitos de probabilidade, a construção e a interpretação de modelos preditivos seriam impossíveis. Modelos de classificação, por exemplo, frequentemente produzem probabilidades (a probabilidade de um email ser spam é 95%). Modelos de regressão, embora prevejam um valor contínuo, são construídos sob a premissa de distribuições de probabilidade dos erros. A probabilidade não é apenas uma teoria abstrata; é a espinha dorsal que sustenta a capacidade de um sistema de Machine Learning de fazer previsões significativas.

Nesta seção, vamos explorar algumas das distribuições de probabilidade mais comuns, que servem como blocos de construção para muitos modelos estatísticos e de Machine Learning. Entender como os dados se distribuem nos ajuda a escolher o modelo certo, a interpretar seus resultados e a quantificar a confiança em nossas previsões.

# Distribuição Normal: A Curva Mais Famosa

📌 **Também conhecida como:** Distribuição Gaussiana ou "Curva em Sino"

A **distribuição normal**, também conhecida como distribuição gaussiana ou "curva em sino", é talvez a distribuição de probabilidade mais importante e amplamente utilizada em estatística e Machine Learning. Ela descreve um fenômeno onde a maioria dos valores se agrupa em torno da média, e os valores se tornam menos frequentes à medida que se afastam dela, simetricamente em ambas as direções. Pense em muitas características naturais, como a altura das pessoas, a pressão sanguínea ou os erros de medição: todas tendem a seguir uma distribuição normal.

## Por que é tão importante?

Sua importância reside no **Teorema do Limite Central**, que afirma que, sob certas condições, a média de um grande número de amostras independentes de qualquer distribuição tenderá a ser normalmente distribuída. Isso é um pilar para a inferência estatística e para muitos algoritmos de Machine Learning que assumem normalidade nos dados ou nos erros. Por exemplo, a regressão linear clássica assume que os resíduos (os erros do modelo) são normalmente distribuídos.

Uma característica chave da distribuição normal é que ela é completamente definida por apenas dois parâmetros: sua **média ( $\mu$ )** e seu **desvio padrão ( $\sigma$ )**. A média determina o centro da curva, e o desvio padrão determina sua "largura" ou dispersão. Quanto menor o desvio padrão, mais "estreita" e "alta" é a curva, indicando que os dados estão mais concentrados em torno da média.

## Parâmetros-chave

- **Média ( $\mu$ ):** Determina o centro da curva
- **Desvio Padrão ( $\sigma$ ):** Determina a largura/dispersão

# Distribuições Binomial e Poisson: Contando Sucessos e Ocorrências

Enquanto a distribuição normal lida com variáveis contínuas, outras distribuições são essenciais para variáveis discretas. A **distribuição binomial** é usada para modelar o número de "sucessos" em uma sequência fixa de  $n$  tentativas independentes, onde cada tentativa tem apenas dois resultados possíveis (sucesso ou fracasso) e a probabilidade de sucesso ( $p$ ) é constante em cada tentativa. Pense em lançar uma moeda 10 vezes e contar o número de caras, ou em testar 50 produtos e contar quantos são defeituosos.



## Distribuição Binomial

**Uso:** Número de sucessos em  $n$  tentativas fixas

**Parâmetros:**  $n$  (tentativas) e  $p$  (probabilidade de sucesso)

**Exemplo:** Número de clientes que clicam em um anúncio de 100 visualizações



## Distribuição de Poisson

**Uso:** Número de eventos em um intervalo fixo

**Parâmetro:**  $\lambda$  (taxa média de ocorrência)

**Exemplo:** Número de e-mails de spam recebidos por hora

Já a **distribuição de Poisson** é ideal para modelar o número de eventos que ocorrem em um intervalo fixo de tempo ou espaço, quando esses eventos acontecem com uma taxa média conhecida e independentemente uns dos outros. Exemplos incluem o número de chamadas recebidas por um call center em uma hora, o número de acidentes em uma rodovia por mês, ou o número de clientes que chegam a uma loja em um dia. Ela é caracterizada por um único parâmetro,  $\lambda$  (lambda), que representa a taxa média de ocorrência dos eventos.

| Conceito | Âmbito/Aplicação                           | Base/Origem   | Exemplo  |
|----------|--|---|--|
| Normal   | Variáveis contínuas, fenômenos naturais    | Média ( $\mu$ ) e Desvio Padrão ( $\sigma$ )                    | Altura de uma população  |
| Binomial | Número de sucessos em $n$ tentativas fixas | Número de tentativas ( $n$ ) e Probabilidade de sucesso ( $p$ ) | Número de clientes que clicam em um anúncio de 100 visualizações |
| Poisson  | Número de eventos em um intervalo fixo     | Taxa média de ocorrência ( $\lambda$ )                          | Número de e-mails de spam recebidos por hora                     |

**XAI e Distribuições:** A escolha da distribuição de probabilidade correta é um passo fundamental na construção de modelos preditivos, especialmente em tarefas de classificação (binomial) ou contagem de eventos (Poisson). A Inteligência Artificial Explicável (XAI) se beneficia dessa compreensão, pois ao entender a distribuição subjacente dos dados, podemos justificar melhor as premissas e os resultados dos modelos, tornando-os mais transparentes e confiáveis.

# Testes de Hipóteses e o Enigmático p-valor: Tomando Decisões com Dados

"Testes de hipóteses são como um julgamento: você precisa de evidências para provar ou refutar uma acusação."

No mundo da ciência de dados e da pesquisa, muitas vezes não estamos apenas descrevendo o que vemos, mas tentando tirar conclusões sobre uma população maior a partir de uma amostra limitada. Queremos saber se uma nova estratégia de marketing realmente aumentou as vendas, se um novo medicamento é mais eficaz que o placebo, ou se uma característica específica realmente influencia o desempenho de um modelo. É aqui que os **testes de hipóteses** entram em cena, oferecendo uma estrutura formal para tomar decisões baseadas em evidências estatísticas.



## Pergunta de Pesquisa

Há um efeito real ou é apenas acaso?



## Coleta de Dados

Obter evidências da amostra



## Cálculo Estatístico

Determinar a probabilidade sob H0



## Decisão

Rejeitar ou não rejeitar H0

Pense nos testes de hipóteses como um julgamento. Você tem uma "acusação" (a hipótese que você quer testar) e precisa de evidências para prová-la ou refutá-la. Sem uma metodologia clara para avaliar essas evidências, nossas conclusões seriam baseadas em intuição ou viés, o que é perigoso em qualquer campo que exija rigor. Os testes de hipóteses nos fornecem um caminho sistemático para avaliar a probabilidade de nossas observações ocorrerem por acaso, ajudando-nos a distinguir entre um efeito real e uma flutuação aleatória.

A capacidade de realizar e interpretar testes de hipóteses é uma habilidade crucial para qualquer profissional de dados. Ela permite validar suposições, comparar grupos, avaliar a significância de variáveis em modelos e, em última instância, tomar decisões mais embasadas e menos propensas a erros.

# A Lógica por Trás dos Testes: Hipótese Nula e Alternativa

## Hipótese Nula (H0)


A base de qualquer teste de hipóteses reside em duas declarações opostas: a **hipótese nula (H0)** e a **hipótese alternativa (H1 ou Ha)**. A hipótese nula é a "posição padrão", a ideia de que não há efeito, não há diferença, ou que qualquer observação é puramente devido ao acaso. É o que assumimos ser verdadeiro até que haja evidências suficientes para o contrário. Por exemplo, H0 pode ser "a nova estratégia de marketing não tem efeito nas vendas".

- Posição padrão
- Não há efeito ou diferença
- Assumida verdadeira até prova em contrário

## Hipótese Alternativa (H1)

A hipótese alternativa é o que você está tentando provar, a declaração de que há um efeito, uma diferença ou uma relação. Ela é o oposto da hipótese nula. No exemplo anterior, H1 seria "a nova estratégia de marketing *umenta* as vendas". O objetivo do teste de hipóteses é usar os dados da amostra para decidir se temos evidências fortes o suficiente para "rejeitar" a hipótese nula em favor da alternativa.

- O que você quer provar
- Há um efeito ou diferença
- Oposto da hipótese nula

 **O Processo:** O processo envolve coletar dados, calcular uma estatística de teste (que mede o quão longe nossos dados estão do que esperaríamos sob a H0) e, em seguida, determinar a probabilidade de observar tal estatística (ou uma mais extrema) se a hipótese nula fosse realmente verdadeira. Essa probabilidade é o famoso p-valor.

# O p-valor Desmistificado: O Que Ele Realmente Significa?

## p-valor $\neq$ probabilidade de H0 ser verdadeira

O **p-valor** é, talvez, um dos conceitos mais mal compreendidos e mal utilizados em estatística. De forma intuitiva, o p-valor é a probabilidade de observar os resultados da sua amostra (ou resultados ainda mais extremos) *se a hipótese nula fosse verdadeira*. Um p-valor pequeno (geralmente menor que 0,05 ou 5%) sugere que seus resultados são improváveis de terem ocorrido por puro acaso, o que nos leva a rejeitar a hipótese nula. Um p-valor grande, por outro lado, indica que seus resultados são bastante plausíveis sob a hipótese nula, e, portanto, não temos evidências suficientes para rejeitá-la.

### O que o p-valor NÃO é

- Não é a probabilidade de H0 ser verdadeira
- Não é a probabilidade de H1 ser verdadeira
- Não é a probabilidade de ter cometido um erro

### O que o p-valor É

- É a probabilidade de observar seus dados (ou mais extremos) se H0 for verdadeira
- É uma medida da força da evidência contra H0
- É um valor entre 0 e 1

É crucial entender que um p-valor não é a probabilidade de a hipótese nula ser verdadeira, nem a probabilidade de a hipótese alternativa ser verdadeira. Ele é uma medida da força da evidência contra a hipótese nula. Se você tem um p-valor de 0,03 para a hipótese de que a nova estratégia de marketing não tem efeito, isso significa que há apenas 3% de chance de você observar um aumento de vendas tão grande (ou maior) se a estratégia realmente não tivesse efeito. Isso é considerado uma evidência forte para rejeitar a hipótese nula e concluir que a estratégia provavelmente funciona.

**Machine Learning e p-valor:** Em Machine Learning, testes de hipóteses são usados para comparar o desempenho de modelos, para selecionar *features* (variáveis) que são estatisticamente significativas para a previsão, ou para validar a eficácia de intervenções. A automação de Machine Learning (AutoML) pode, em alguns casos, incorporar testes de hipóteses para a seleção de modelos ou otimização de hiperparâmetros, embora muitas vezes opere em um nível mais empírico. No entanto, a compreensão humana do p-valor é insubstituível para a interpretabilidade e a validação crítica dos resultados, especialmente em contextos de Inteligência Artificial Explicável (XAI), onde justificar as decisões do modelo é tão importante quanto a própria previsão.

# Consolidação: A Estatística como Sua Aliada na Modelagem Preditiva

Chegamos ao fim da nossa revisão de conceitos essenciais de estatística, mas este é apenas o começo da sua aplicação prática. Vimos que as medidas de tendência central e dispersão são como os olhos e ouvidos para entender o comportamento básico dos seus dados. As distribuições de probabilidade nos dão a linguagem para quantificar a incerteza e modelar fenômenos, enquanto os testes de hipóteses e o p-valor nos fornecem um arcabouço rigoroso para tomar decisões informadas a partir de evidências.

|   |   |
|---|---|
| <b>Análise Exploratória</b><br>Use média, mediana, moda, variância e desvio padrão para entender suas variáveis | <b>Identificação de Distribuições</b><br>Reconheça a distribuição de suas features e variáveis-alvo |
| <b>Seleção de Modelos</b><br>Escolha modelos adequados baseados nas características dos dados                   | <b>Validação Estatística</b><br>Use testes de hipóteses para garantir conclusões válidas            |

## Em prática:

Ao iniciar um projeto de Machine Learning, comece sempre com uma análise exploratória de dados robusta, utilizando média, mediana, moda, variância e desvio padrão para entender suas variáveis. Identifique a distribuição de suas *features* e variáveis-alvo para escolher os modelos mais adequados. E, ao comparar modelos ou avaliar a significância de variáveis, utilize testes de hipóteses para garantir que suas conclusões são estatisticamente válidas, não apenas coincidências.

# Autoavaliação

1

## Questão 1

Qual das seguintes medidas de tendência central é mais robusta à presença de *outliers* em um conjunto de dados salariais?

- a) Média
- b) Mediana
- c) Moda
- d) Desvio Padrão

2

## Questão 2

Um cientista de dados está analisando o tempo de resposta de um servidor e observa que os dados estão muito espalhados. Qual medida estatística ele deve usar para quantificar essa variabilidade na mesma unidade de tempo?

- a) Variância
- b) Média
- c) Desvio Padrão
- d) Moda

3

## Questão 3

A distribuição de probabilidade mais adequada para modelar o número de clientes que entram em uma loja por hora, sabendo-se a taxa média de entrada, é a:

- a) Normal
- b) Binomial
- c) Poisson
- d) Uniforme

4

## Questão 4

Em um teste de hipóteses, um p-valor de 0,01 significa que:

- a) Há 1% de chance de a hipótese nula ser verdadeira.
- b) Há 99% de chance de a hipótese alternativa ser verdadeira.
- c) Há 1% de chance de observar os resultados da amostra (ou mais extremos) se a hipótese nula for verdadeira.
- d) A hipótese nula deve ser aceita.

5

## Questão 5

Explique a importância de entender as medidas de dispersão (variância e desvio padrão) para a fase de pré-processamento de dados em um projeto de Machine Learning.

## Gabarito

1. **b) Mediana**
2. **c) Desvio Padrão**
3. **c) Poisson**
4. c) Há 1% de chance de observar os resultados da amostra (ou mais extremos) se a hipótese nula for verdadeira.

# Próximos Passos e Recursos



## Próxima Aula

Aula 3 – Regressão Linear e Logística: A Base dos Modelos

**Próxima Aula:** Na Aula 3 – Regressão Linear e Logística: A Base dos Modelos, daremos o próximo passo, aplicando muitos dos conceitos estatísticos revisados hoje para construir e interpretar os modelos preditivos mais fundamentais.

## Recursos Adicionais



### Livro Recomendado

**"Practical Statistics for Data Scientists"**

Para aprofundar a conexão entre estatística e ciência de dados.



### Plataforma de Estudo

**Khan Academy - Estatística e Probabilidade**

Para revisar conceitos básicos de forma interativa.



### Artigos Técnicos

**Artigos sobre XAI e AutoML**

Para entender como esses conceitos se integram às tendências atuais.



**NOTA IMPORTANTE:** As informações regulatórias/legais/técnicas desta aula estão atualizadas até 2025. Consulte sempre fontes oficiais para verificar alterações.