

Aula 8 – Seleção e Filtragem de Dados

Imagine-se diante de uma montanha de informações, um verdadeiro oceano de dados. Pode ser uma planilha gigantesca de vendas, um banco de dados de pacientes ou até mesmo os resultados de uma pesquisa de opinião. Em meio a tanta informação, como você encontra exatamente o que precisa para tomar uma decisão inteligente ou para responder a uma pergunta crucial? É aqui que a **seleção e filtragem de dados** entram em cena, transformando o caos em clareza.

Introdução – O Poder de Focar no Essencial

Nosso Objetivo Principal

Que você seja capaz de navegar por grandes conjuntos de dados com confiança, extraíndo informações valiosas de forma eficiente e reproduzível.

Ferramentas Poderosas

Vamos explorar as ferramentas mais poderosas do Python, como a biblioteca Pandas, que são o padrão da indústria para análise de dados.

Transformação

Prepare-se para transformar a maneira como você interage com as informações, abrindo portas para análises mais profundas e decisões mais assertivas.

Nesta aula, embarcaremos em uma jornada para dominar as técnicas que permitem "conversar" com seus dados, pedindo a eles que mostrem apenas o que é relevante. Você aprenderá a isolar colunas específicas, a escolher linhas que atendam a critérios precisos e a combinar essas operações para desvendar padrões ocultos. Ao final, você não apenas saberá como manipular dados, mas entenderá o "porquê" por trás de cada escolha, tornando-se um analista mais estratégico e eficaz.

- 📄 Conectando com o que você já conhece, pense em como você organiza seus e-mails ou arquivos no computador. Você não lê todos os e-mails para encontrar um específico, certo? Você usa filtros! Da mesma forma, em planilhas, você aplica filtros para ver apenas os dados de um determinado período ou cliente. Agora, vamos levar essa lógica para o universo da programação, ganhando um poder e uma flexibilidade inimagináveis.

O Desafio da Montanha de Dados e a Necessidade de Foco

No mundo atual, somos constantemente bombardeados por dados. Empresas coletam informações sobre seus clientes, governos registram dados demográficos e econômicos, e cientistas geram volumes massivos de resultados de experimentos. Essa abundância, embora valiosa, pode ser esmagadora. Tentar analisar todos os dados de uma vez é como tentar beber água de uma mangueira de incêndio: ineficiente e improdutivo.

01

O Problema Central

Nem toda informação é igualmente relevante para a pergunta que você está tentando responder.

02

A Solução

A chave para uma análise eficaz reside na capacidade de **focar**.

03

A Ferramenta

Seleção e filtragem de dados se tornam habilidades indispensáveis.

É aqui que a seleção e a filtragem de dados se tornam habilidades indispensáveis. Elas permitem que você corte o "ruído", isolando apenas as informações que realmente importam para sua análise. Pense nisso como um garimpeiro que, em vez de carregar toda a terra do rio, desenvolve técnicas para separar o ouro das pedras e da areia.

No contexto da análise de dados com Python, essa "terra" é frequentemente representada por um **DataFrame** da biblioteca Pandas. Um DataFrame é como uma tabela, com linhas e colunas, onde cada coluna tem um nome e cada linha representa uma observação. Nossa missão, então, é aprender a "garimpar" esse DataFrame, extraíndo o "ouro" que nos levará a insights valiosos.

Selecionando Colunas – O Foco da Sua Análise

Quando você olha para uma planilha, seus olhos naturalmente buscam as colunas que contêm as informações mais importantes para sua tarefa. Se você está interessado no nome do cliente e no valor da compra, você ignora a coluna de endereço ou de data de nascimento, certo? No Pandas, fazemos exatamente isso, mas de forma programática e muito mais eficiente.

Por que Selecionar Colunas?

- Remove o excesso de informação
- Direciona sua atenção para o que realmente importa
- Torna seu código mais limpo
- Otimiza o desempenho com grandes volumes de dados

Analogia

É como escolher quais ingredientes você vai usar para uma receita: você não pega tudo da geladeira, mas apenas o que é necessário para aquele prato específico.

```
# Selecionando uma única coluna
coluna_produto = df_vendas['Produto']

# Selecionando múltiplas colunas
dados_essenciais = df_vendas[['Produto', 'Valor', 'Cliente']]
```

A maneira mais comum e intuitiva de selecionar colunas em um DataFrame Pandas é usando colchetes []. Se você quer uma única coluna, basta passar o nome dela como uma string. Se precisar de múltiplas colunas, você passa uma lista de strings com os nomes das colunas desejadas.

- ❏ No ambiente profissional, isso significa poder rapidamente gerar relatórios focados em métricas chave, como o desempenho de vendas por produto ou a receita total por cliente, sem se perder em detalhes irrelevantes.

Selecionando Linhas com .loc[] – Onde a Magia Acontece

Depois de decidir quais colunas você quer ver, o próximo passo é escolher quais linhas são relevantes. Afinal, mesmo com as colunas certas, você pode ter milhares ou milhões de registros, e nem todos eles se encaixam no seu critério de análise. É como ter um mapa com todas as ruas, mas você só quer ver as ruas que estão dentro de um bairro específico.

O que é .loc[]?

O método .loc[] no Pandas é a sua ferramenta para selecionar linhas e colunas com base em seus **rótulos** (ou nomes).

Flexibilidade

Você pode usá-lo para selecionar uma única linha, um intervalo de linhas, ou até mesmo linhas que atendam a uma condição lógica.

Analogia

Pense nele como um sistema de coordenadas onde você especifica o nome da linha e o nome da coluna que deseja acessar.

```
# Selecionando a linha com o rótulo (ID) 'Venda_001'
venda_especifica = df_vendas.loc['Venda_001']

# Selecionando um intervalo de linhas por rótulo
# Note que o .loc[] inclui o final do intervalo
intervalo_vendas = df_vendas.loc['Venda_001':'Venda_005']

# Selecionando linhas que atendem a uma condição e colunas específicas
vendas_grandes_sp = df_vendas.loc[df_vendas['Valor'] > 1000, ['Produto', 'Valor', 'Cliente']]
```

No dia a dia de um analista, o .loc[] é fundamental para tarefas como: encontrar todos os registros de um cliente específico, analisar transações realizadas em um determinado mês (se o índice for data), ou extrair dados de uma categoria de produto específica. Ele oferece uma maneira poderosa e legível de focar sua análise em subconjuntos de dados relevantes.

Selecionando Linhas por Posição com `.iloc[]` – A Precisão Numérica

Enquanto o `.loc[]` trabalha com os rótulos (nomes) das linhas e colunas, o `.iloc[]` (de "integer location") opera com base nas **posições inteiras** dos elementos. Pense nele como um sistema de coordenadas cartesianas, onde você especifica a linha e a coluna pelo seu número de ordem, começando do zero. É como um elevador que vai para o andar número 5, independentemente do nome que esse andar possa ter.

1

`.loc[]` - Por Rótulos

Trabalha com nomes e rótulos das linhas e colunas

2

`.iloc[]` - Por Posição

Opera com posições inteiras, começando do zero

```
# Selecionando a primeira linha (índice 0)
primeira_venda = df_vendas.iloc[0]

# Selecionando as primeiras 5 linhas (do índice 0 ao 4)
# Note que o 5 é exclusivo, ou seja, vai até o índice 4
cinco_primeiras_vendas = df_vendas.iloc[0:5]

# Selecionando linhas e colunas específicas por seus índices
# Por exemplo, a 1ª, 3ª e 5ª linha, e a 2ª e 4ª coluna
dados_selecionados_por_posicao = df_vendas.iloc[[0, 2, 4], [1, 3]]
```

- ❏ O `.iloc[]` é extremamente útil para tarefas como inspecionar os primeiros ou últimos registros de um conjunto de dados para entender sua estrutura, ou para extrair amostras de dados baseadas em sua ordem de ocorrência.

.loc[] vs. .iloc[] – Escolhendo a Ferramenta Certa

A essa altura, você deve estar se perguntando: "Quando devo usar .loc[] e quando devo usar .iloc[]?". Essa é uma pergunta excelente e fundamental para dominar a manipulação de dados com Pandas. A escolha entre um e outro depende diretamente de como você quer referenciar seus dados: pelos seus nomes (rótulos) ou por suas posições (índices numéricos).

Característica	.loc[]	.iloc[]
Base de Seleção	Rótulos (nomes) de linhas e colunas	Posições inteiras (índices) de linhas e colunas
Tipo de Índice	Pode ser qualquer tipo (string, int, datetime)	Apenas inteiros (0, 1, 2, ...)
Fatiamento	Inclui o valor final (ex: 0:5 inclui 5)	Exclui o valor final (ex: 0:5 vai até 4)
Uso Comum	Seleção por condições, por rótulos específicos	Seleção dos primeiros/últimos N, por posição
Analogia	Buscar por nome/título	Buscar por posição/número na fila

Pense na diferença como buscar um livro em uma biblioteca. Se você sabe o título do livro e o nome da estante onde ele está, você usa o .loc[] para ir diretamente a essa "coordenada nomeada". Mas se você sabe que o livro que você quer é o terceiro livro da quinta prateleira, você usa o .iloc[] para ir àquela "posição numérica" exata.

Compreender essa distinção é crucial para escrever código Pandas robusto e claro. No ambiente de trabalho, escolher a ferramenta certa evita erros sutis e torna sua análise mais intuitiva para quem a lê.

Filtragem Baseada em Condições Lógicas

– O Poder do "E se..."

Até agora, aprendemos a selecionar colunas e linhas por seus nomes ou posições. Mas e se você quiser encontrar todas as vendas que foram maiores que R\$ 500, ou todos os clientes que moram em São Paulo e têm mais de 30 anos? É aqui que a **filtragem baseada em condições lógicas** se torna a verdadeira estrela da análise de dados.



O Superpoder

A filtragem por condição é como ter um superpoder de "perguntar" aos seus dados. Você não está mais apenas pegando pedaços, mas sim definindo regras para que os dados se revelem.



Máscara Booleana

No Pandas, essa "pergunta" é feita criando uma máscara booleana - uma Série composta apenas por valores True ou False.



Resultado Focado

Quando você passa essa máscara para o DataFrame, ele magicamente retorna apenas as linhas onde a máscara é True.

```
# Criando uma máscara booleana para vendas com valor maior que 500
condicao_valor_alto = df_vendas['Valor'] > 500
print(condicao_valor_alto.head()) # Isso mostraria True/False para cada linha

# Usando a máscara para filtrar o DataFrame
vendas_altas = df_vendas[condicao_valor_alto]

# Ou, de forma mais concisa, diretamente:
vendas_altas_direto = df_vendas[df_vendas['Valor'] > 500]
```

Essa técnica é a espinha dorsal de quase toda análise exploratória de dados. Ela permite que você isole grupos específicos para estudo, identifique anomalias, ou prepare dados para visualizações que respondam a perguntas muito específicas. No mundo real, isso significa poder identificar rapidamente clientes de alto valor, produtos com baixo desempenho ou regiões com maior potencial de crescimento.

Combinando Múltiplas Condições – O "E" (&) e o "OU" (|)

Raramente uma única condição é suficiente para isolar o subconjunto de dados que você realmente precisa. Na vida real, nossas perguntas são mais complexas: "Quero ver os clientes que moram em São Paulo *E* que gastaram mais de R\$ 1000", ou "Preciso de todos os produtos que estão em promoção *OU* que têm estoque baixo". Para isso, precisamos combinar condições.

Operador & (E/AND)

Os dados precisam satisfazer **todas** as condições para o resultado ser True

Operador | (OU/OR)

Os dados precisam satisfazer **pelo menos uma** das condições para o resultado ser True

Parênteses ()

É **crucial** envolver cada condição individual entre parênteses para garantir a precedência correta

```
# Filtrando vendas de produtos com valor acima de 500 E da região 'Sudeste'  
vendas_sudeste_altas = df_vendas[(df_vendas['Valor'] > 500) & (df_vendas['Regiao'] == 'Sudeste')]
```

```
# Filtrando vendas de produtos 'Notebook' OU 'Smartphone'  
vendas_eletronicos = df_vendas[(df_vendas['Produto'] == 'Notebook') | (df_vendas['Produto'] == 'Smartphone')]
```

```
# Combinando três condições: Valor > 1000 E Regiao == 'Sul' OU Produto == 'TV'  
# Note a importância dos parênteses para agrupar as operações  
condicao_complexa = df_vendas[((df_vendas['Valor'] > 1000) & (df_vendas['Regiao'] == 'Sul')) | (df_vendas['Produto'] == 'TV')]
```

Pense nisso como um sistema de segurança onde você precisa de duas chaves para abrir uma porta (&), ou onde qualquer uma das duas chaves já é suficiente (|). A lógica é a mesma: os dados precisam satisfazer *todas* as condições para o & ser True, ou *pelo menos uma* das condições para o | ser True.

Dominar a combinação de condições é o que eleva sua capacidade de análise a um novo patamar. Você pode segmentar seus dados de forma incrivelmente granular, identificando nichos de mercado, problemas específicos de produtos ou oportunidades de vendas que seriam impossíveis de encontrar sem essa flexibilidade.

A Elegância do Método `.query()`

Embora a filtragem com máscaras booleanas usando `[]` seja poderosa, a sintaxe pode se tornar um pouco densa e difícil de ler, especialmente quando você combina muitas condições. Imagine ter que escrever `df[(df['coluna1'] > 10) & (df['coluna2'] == 'A') | (df['coluna3'].isin(['X', 'Y']))]` – os parênteses e as repetições de `df['coluna']` podem confundir.

Sintaxe Tradicional

```
df[(df['Valor'] > 500) &
   (df['Regiao'] == 'Sudeste')]
```

Com `.query()`

```
df.query('Valor > 500 and
         Regiao == "Sudeste"')
```

Para tornar a filtragem mais legível e "Pythonic", o Pandas oferece o método `.query()`. Ele permite que você escreva suas condições de filtragem como uma string, usando nomes de colunas diretamente, como se estivesse escrevendo uma consulta em SQL. É como passar de uma linguagem de programação mais "robótica" para uma linguagem mais próxima do inglês, tornando seu código mais intuitivo.

```
# Filtrando vendas com valor acima de 500 E da região 'Sudeste'
vendas_sudeste_altas_query = df_vendas.query('Valor > 500 and Regiao == "Sudeste"')

# Filtrando vendas de produtos 'Notebook' OU 'Smartphone'
vendas_eletronicos_query = df_vendas.query('Produto == "Notebook" or Produto == "Smartphone"')

# Combinando três condições: Valor > 1000 E Regiao == 'Sul' OU Produto == 'TV'
condicao_complexa_query = df_vendas.query('(Valor > 1000 and Regiao == "Sul") or Produto == "TV"')
```

- 📌 Note como a sintaxe se torna mais limpa e direta. Você usa **and** e **or** em vez de `&` e `|`, e os nomes das colunas são usados diretamente. Para strings, você pode usar aspas simples ou duplas. Se você precisar referenciar uma variável Python dentro da sua consulta, basta prefixá-la com `@`.

O `.query()` é uma ferramenta de produtividade que melhora a legibilidade do seu código, o que é inestimável em projetos colaborativos ou quando você revisita seu próprio código meses depois. Ele não substitui a filtragem com máscaras booleanas, mas oferece uma alternativa elegante para cenários específicos.

Seleção e Filtragem na Prática – Um Estudo de Caso Simples

Agora que exploramos as ferramentas individualmente, vamos juntá-las em um cenário prático. Imagine que você é um analista de dados em uma startup de e-commerce e precisa entender o desempenho de vendas de produtos eletrônicos na região Sudeste, especificamente para vendas acima de R\$ 700.

01

Carregar os Dados

Primeiro, vamos simular a criação de um DataFrame para nosso exemplo

02

Selecionar Colunas Relevantes

Para nossa análise, 'Produto', 'Valor', 'Regiao' e 'Categoria' são as mais importantes

03

Filtrar por Múltiplas Condições

Aplicar as condições: 'Categoria' é 'Eletrônicos', 'Regiao' é 'Sudeste' e 'Valor' é maior que 700

```
import pandas as pd

dados = {
    'Produto': ['Notebook', 'Smartphone', 'TV', 'Fone', 'Notebook', 'Smartphone', 'TV', 'Fone'],
    'Valor': [1200, 800, 1500, 150, 1300, 900, 1600, 180],
    'Cliente': ['Ana', 'Bruno', 'Carlos', 'Diana', 'Eduardo', 'Fernanda', 'Gustavo', 'Helena'],
    'Data': pd.to_datetime(['2024-01-10', '2024-01-12', '2024-01-15', '2024-01-18',
                           '2024-02-01', '2024-02-05', '2024-02-10', '2024-02-14']),
    'Regiao': ['Sudeste', 'Sul', 'Nordeste', 'Sudeste', 'Sudeste', 'Sul', 'Nordeste', 'Sudeste'],
    'Categoria': ['Eletrônicos', 'Eletrônicos', 'Eletrônicos', 'Acessórios',
                 'Eletrônicos', 'Eletrônicos', 'Eletrônicos', 'Acessórios']
}

df_vendas = pd.DataFrame(dados)
```

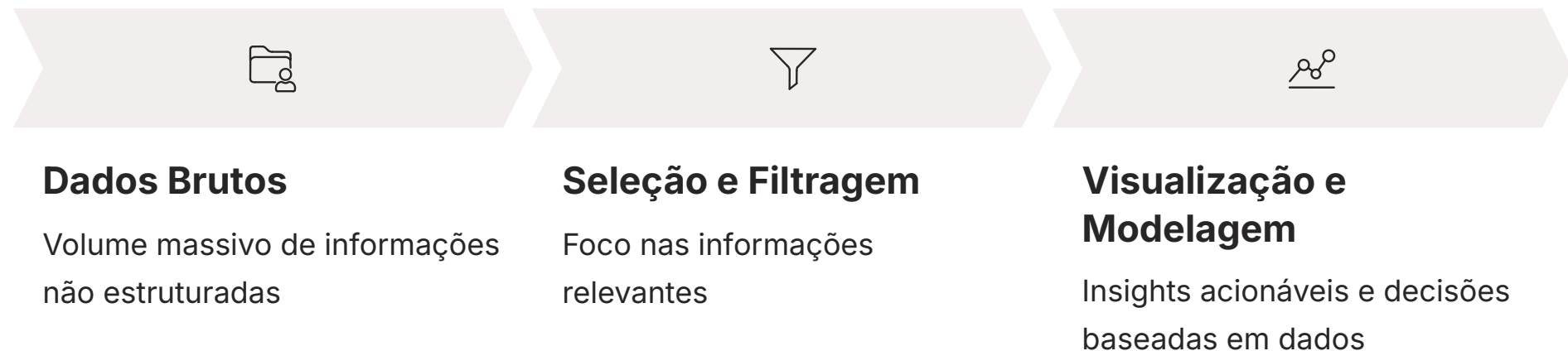
```
# Passo 2: Selecionar Colunas Relevantes
df_foco = df_vendas[['Produto', 'Valor', 'Regiao', 'Categoria']]

# Passo 3: Filtrar por Múltiplas Condições
vendas_eletronicos_sudeste_alto_valor = df_foco.query(
    'Categoria == "Eletrônicos" and Regiao == "Sudeste" and Valor > 700'
)
```

Este fluxo de trabalho, que combina seleção de colunas e filtragem de linhas, é a base para qualquer análise de dados. Ele permite que você refine seu foco, transformando um grande conjunto de dados em um subconjunto gerenciável e altamente relevante para suas perguntas de negócio.

O Impacto da Seleção e Filtragem na Análise Exploratória

A seleção e filtragem de dados não são apenas etapas isoladas; elas são o coração da **Análise Exploratória de Dados (AED)**. Pense nelas como as lentes que você usa para enxergar melhor. Sem elas, seus dados seriam apenas um borrão indistinto. Com elas, você pode focar em detalhes específicos, identificar padrões, detectar anomalias e, finalmente, formular hipóteses sobre o que os dados estão tentando lhe dizer.



Quando você filtra dados, você está essencialmente criando um subconjunto que é mais fácil de visualizar e modelar. Por exemplo, se você quer criar um gráfico de barras mostrando as vendas por categoria de produto, faz sentido primeiro filtrar apenas as colunas 'Produto' e 'Valor'. Se você quer analisar a distribuição de idade dos seus clientes mais fiéis, você primeiro filtra os clientes fiéis e depois seleciona a coluna 'Idade'.

A analogia aqui é a de um chef de cozinha. Antes de começar a cozinhar, ele não joga todos os ingredientes da despensa na panela. Ele seleciona os ingredientes certos para a receita (seleção de colunas) e os prepara (lava, corta, filtra impurezas – filtragem de linhas). Só então ele começa a cozinhar (análise e modelagem).

Essa preparação é crucial para as próximas etapas da análise. Ferramentas de visualização como Matplotlib, Seaborn e Plotly (que são o padrão da indústria e serão abordadas no curso) funcionam muito melhor com dados limpos e focados. Tentar plotar um DataFrame inteiro com dezenas de colunas e milhões de linhas pode ser ineficiente ou até mesmo impossível de interpretar.

Reprodutibilidade e Jupyter Notebooks

No mundo da análise de dados, não basta apenas chegar a uma conclusão; é fundamental que o processo que levou a essa conclusão seja **reproduzível**. Isso significa que qualquer pessoa, ou até mesmo você mesmo no futuro, deve ser capaz de executar os mesmos passos e obter os mesmos resultados. A seleção e filtragem de dados, quando feitas em um ambiente como o Jupyter Notebook, são pilares da reprodutibilidade.

O que é um Jupyter Notebook?

Um ambiente interativo onde você pode combinar código (Python), texto explicativo (Markdown), equações e visualizações em um único documento.

Registro Claro

Ao realizar suas operações de seleção e filtragem dentro de células de código, você cria um registro claro e executável de cada passo da sua análise.

Colaboração Eficiente

Em vez de explicar verbalmente ou com capturas de tela estáticas, você pode simplesmente compartilhar seu Jupyter Notebook.

Imagine que você está trabalhando em um projeto e precisa mostrar a um colega como você chegou a um determinado subconjunto de dados. Em vez de explicar verbalmente ou com capturas de tela estáticas, você pode simplesmente compartilhar seu Jupyter Notebook. Seu colega pode executar as células, ver os dados brutos, as condições de filtragem aplicadas e os resultados intermediários e finais. Isso não só economiza tempo, mas também aumenta a confiança na sua análise.

- ❏ A capacidade de reproduzir uma análise é crucial para a validação, para a colaboração em equipe e para a manutenção de projetos de dados a longo prazo. Ao usar Pandas para selecionar e filtrar seus dados em um Jupyter Notebook, você está construindo um "laboratório" de dados transparente e verificável, um padrão de excelência no mercado de trabalho atual.

Storytelling com Dados – O Que a Filtragem Revela

A análise de dados não se resume a números e gráficos; ela é sobre contar uma **história**. E, assim como um bom contador de histórias seleciona cuidadosamente os detalhes e os eventos que compõem a narrativa, um bom analista de dados usa a seleção e filtragem para revelar a trama principal que os dados escondem.

01

Dados Brutos = Palavras Embaralhadas

Dados brutos são como um livro com todas as palavras embaralhadas

02

Seleção e Filtragem = Organização

São as ferramentas que permitem organizar essas palavras em frases, parágrafos e capítulos coerentes

03

Resultado = História Clara

Uma narrativa clara e convincente, essencial para comunicar seus achados

Por exemplo, se você tem dados de feedback de clientes, filtrar apenas os comentários negativos de clientes que cancelaram o serviço pode revelar uma história clara sobre um ponto de dor específico no seu produto ou atendimento. Sem essa filtragem, esses insights estariam perdidos em meio a milhares de outros comentários.

A capacidade de filtrar e focar em subconjuntos específicos de dados permite que você construa argumentos mais fortes e apresente descobertas mais impactantes. Em vez de dizer "as vendas variaram", você pode dizer "as vendas de produtos premium para clientes novos na região X caíram 15% no último trimestre devido a [causa identificada após filtragem]".

Pense em um detetive que, diante de um mar de pistas, filtra as irrelevantes e foca nas que realmente apontam para o culpado. Ele não apresenta todas as pistas, mas sim a sequência lógica que o levou à conclusão. Da mesma forma, a seleção e filtragem são suas ferramentas para destilar a complexidade dos dados em uma narrativa clara e convincente, essencial para comunicar seus achados a stakeholders e tomar decisões baseadas em evidências.

Desafios Comuns e Dicas de Ouro

Ao trabalhar com seleção e filtragem, alguns desafios são comuns, especialmente para quem está começando. Estar ciente deles pode economizar horas de depuração e frustração.

Erro Comum: Parênteses

Esquecer os **parênteses** ao combinar múltiplas condições com & ou |. Sem eles, o Python pode interpretar a expressão de forma diferente, levando a resultados incorretos.

Erro Comum: Tipos de Dados

Verificação dos **tipos de dados**. Se você tentar filtrar uma coluna numérica usando uma string, o Pandas pode não encontrar correspondências.

Erro Comum: Fatiamento

Lembre-se da diferença: `.loc[]` inclui o ponto final do intervalo, enquanto `.iloc[]` o exclui.



Dica de Ouro 1

Comece Pequeno: Se a filtragem é complexa, construa-a passo a passo. Crie uma máscara booleana por vez e visualize-a antes de combiná-las.



Dica de Ouro 2

Use `.copy()`: Ao criar um subconjunto de um DataFrame para modificação, use `.copy()` para evitar o "SettingWithCopyWarning".



Dica de Ouro 3

Documente: Use comentários no seu código ou células de Markdown no Jupyter Notebook para explicar a lógica por trás de suas seleções e filtragens.

Com essas técnicas de seleção e filtragem, você está pronto para o próximo passo: lidar com os dados que *não* estão lá. Na próxima aula, mergulharemos no fascinante mundo do **Tratamento de Dados Ausentes (Missing Values)**, um desafio comum e crucial na análise de dados.

Consolidação e Próximos Passos

Chegamos ao fim de uma jornada essencial no universo da análise de dados. Nesta aula, desvendamos o poder da seleção e filtragem, aprendendo a focar no que realmente importa em meio a um mar de informações. Vimos como as ferramentas do Pandas, como `[]`, `.loc[]`, `.iloc[]`, e o elegante `.query()`, nos permitem extrair insights precisos, seja selecionando colunas específicas, linhas por rótulo ou posição, ou aplicando condições lógicas complexas.

Seleção de Colunas

Sempre comece sua análise exploratória selecionando as colunas mais relevantes

Método `.query()`

Considere usar `.query()` para tornar suas condições de filtragem mais legíveis



`.loc[]` para Rótulos

Use `.loc[]` quando precisar de linhas e colunas por seus nomes/rótulos

`.iloc[]` para Posições

Opte por `.iloc[]` para acesso posicional, como os primeiros ou últimos registros

Máscaras Booleanas

Domine as máscaras booleanas e a combinação de condições com `&` e `|`

Compreendemos que a seleção e filtragem não são apenas operações técnicas, mas habilidades estratégicas que preparam seus dados para visualizações impactantes, modelos preditivos eficazes e, acima de tudo, para contar histórias convincentes. A capacidade de segmentar e refinar seus dados é o que transforma um conjunto bruto de informações em uma fonte de conhecimento acionável.

Autoavaliação

Questões de Múltipla Escolha

- Qual método do Pandas é mais adequado para selecionar linhas com base em seus rótulos (nomes) e colunas específicas por seus nomes?**
 - a) `df[]`
 - b) `df.iloc[]`
 - c) `df.loc[]`
 - d) `df.query()`
- Para filtrar um DataFrame `df` e obter apenas as linhas onde a coluna 'Idade' é maior que 25 E a coluna 'Cidade' é igual a 'Rio de Janeiro', qual das seguintes sintaxes está **correta**?**
 - a) `df[df['Idade'] > 25 & df['Cidade'] == 'Rio de Janeiro']`
 - b) `df[(df['Idade'] > 25) & (df['Cidade'] == 'Rio de Janeiro')]`
 - c) `df.query('Idade > 25' and 'Cidade == "Rio de Janeiro")`
 - d) `df.iloc[df['Idade'] > 25, df['Cidade'] == 'Rio de Janeiro']`
- Se você tem um DataFrame `df` e quer acessar a segunda linha e a terceira coluna usando seus índices numéricos (começando do zero), qual seria a sintaxe correta?**
 - a) `df.loc[1, 2]`
 - b) `df.iloc[1, 2]`
 - c) `df[1][2]`
 - d) `df.query('index == 1 and column == 2')`
- O que o método `.query()` oferece de principal vantagem em relação à filtragem tradicional com máscaras booleanas?**
 - a) Maior velocidade de execução para grandes DataFrames.
 - b) Capacidade de filtrar por múltiplos DataFrames simultaneamente.
 - c) Sintaxe mais legível e próxima da linguagem natural para condições.
 - d) Permite a seleção de colunas e linhas por rótulos e posições ao mesmo tempo.

Questão Discursiva

Explique, com suas palavras, a importância da reprodutibilidade na análise de dados e como a seleção e filtragem, quando realizadas em um Jupyter Notebook, contribuem para esse princípio.

Gabarito e Próximos Passos

Gabarito:

1. c) `df.loc[]`
2. b) `df[(df['Idade'] > 25) & (df['Cidade'] == 'Rio de Janeiro')]`
3. b) `df.iloc[1, 2]`
4. c) Sintaxe mais legível e próxima da linguagem natural para condições.

Resposta Sugerida - Questão Discursiva:

A reprodutibilidade na análise de dados é crucial porque garante que os resultados de uma análise possam ser verificados e replicados por qualquer pessoa, a qualquer momento. Isso aumenta a confiança nos achados e facilita a colaboração. A seleção e filtragem, quando realizadas em um Jupyter Notebook, contribuem para isso ao registrar cada passo da manipulação dos dados em células de código executáveis. Isso cria um "caminho" transparente e auditável, permitindo que outros (ou você mesmo no futuro) sigam exatamente a mesma lógica e obtenham os mesmos subconjuntos de dados, validando assim a análise.

Próxima Aula

Aula 9 – Tratamento de Dados Ausentes (Missing Values)

Recursos Adicionais

- Documentação Oficial do Pandas
- Livro "Python for Data Analysis" (Wes McKinney)
- Artigos no Medium/Towards Data Science

Nota Importante

As informações técnicas desta aula estão atualizadas até 2025. Consulte sempre a documentação oficial das bibliotecas para verificar alterações ou novas funcionalidades.