

Programação para Ciência de Dados

Introdução ao Pandas - Series e DataFrames

Arthur Casals

21 de Outubro de 2025

- ▶ Primeira entrega quinzenal: entregue!
- ▶ Recursos adicionais: notebooks e alternativas a NumPy/Pandas
- ▶ Presença: Moodle
- ▶ Interação via Zoom

Agenda

- ▶ Estruturas Pandas: Series e DataFrames
- ▶ Visualização e Inspeção de Dados
- ▶ Seleção e Indexação
- ▶ Filtragem e Ordenação

Recapitulação Aulas 03-04

O que já sabemos sobre NumPy:

- ▶ Arrays multidimensionais eficientes
- ▶ Operações vetorizadas e broadcasting
- ▶ Indexação avançada (boolean, fancy)
- ▶ Álgebra linear (dot, solve, inv, eig)
- ▶ Manipulação de shape (reshape, transpose)
- ▶ Concatenação e stacking
- ▶ Persistência (save, load, npz)

Hoje vamos aprender:

- ▶ Pandas: biblioteca para análise de dados
- ▶ Series e DataFrames: estruturas de alto nível
- ▶ Manipulação de dados tabulares

Bloco 1

Estruturas Pandas: Series e DataFrames

O que é Pandas?

Pandas (Python Data Analysis Library):

- ▶ Biblioteca para análise e manipulação de dados
- ▶ Construída sobre NumPy
- ▶ Estruturas de dados de alto nível
- ▶ Ferramentas poderosas para trabalhar com dados tabulares
- ▶ Integração com Matplotlib para visualização

Principais características:

- ▶ **Series:** Array 1D rotulado
- ▶ **DataFrame:** Tabela 2D (linhas e colunas rotuladas)
- ▶ **Indexação sofisticada:** Labels em vez de apenas posições
- ▶ **Operações de grupo:** GroupBy poderoso
- ▶ **Merge/Join:** Combinar datasets facilmente
- ▶ **Time Series:** Funcionalidades temporais avançadas

Por que Pandas?

NumPy:

- ▶ Arrays numéricos homogêneos
- ▶ Performance máxima
- ▶ Operações matemáticas
- ▶ Base para computação científica
- ▶ Indexação por posição

Pandas:

- ▶ Dados heterogêneos (tipos mistos)
- ▶ Operações de análise de dados
- ▶ Limpeza e transformação
- ▶ Dados tabulares (como Excel/SQL)
- ▶ Indexação por rótulos

💡 Nota Importante

Pandas é ideal para análise exploratória de dados (EDA)

Quando usar cada um?

Use NumPy quando:

- ▶ Precisar de máxima performance
- ▶ Trabalhar com arrays numéricos puros
- ▶ Fazer álgebra linear
- ▶ Implementar algoritmos matemáticos
- ▶ Processar imagens ou sinais

Use Pandas quando:

- ▶ Trabalhar com dados tabulares (CSV, Excel, SQL)
- ▶ Precisar de rótulos nas linhas/colunas
- ▶ Fazer análise exploratória
- ▶ Limpar e transformar dados
- ▶ Agrupar, agregar e pivotar dados
- ▶ Trabalhar com séries temporais

Instalação e Importação

Instalação:

</> Bash

```
1 # Google Colab: ja vem instalado
2 # Local (via pip):
3 pip install pandas
4
5 # Verificar versao:
6 pip show pandas
7
```

Instalação e Importação

Importação (convenção universal):

</> Python

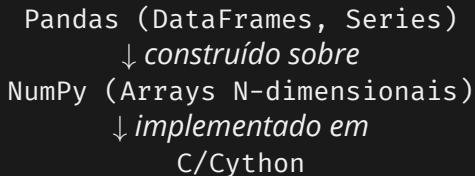
```
1 import pandas as pd
2 import numpy as np # Frequentemente usado junto
3
4 # Verificar versao
5 print(pd.__version__) # 2.0.3 (ou similar)
6
```

💡 Nota Importante

Sempre use `import pandas as pd` - é a convenção padrão

Pandas: Camada sobre NumPy

Arquitetura:



Implicações:

- ▶ DataFrames internamente usam arrays NumPy
- ▶ Pandas herda a performance do NumPy
- ▶ Conversão entre Pandas e NumPy é trivial
- ▶ Muitas operações vetorizadas disponíveis
- ▶ Integração perfeita com ecossistema científico

Series: Array 1D Rotulado

O que é uma Series?

- ▶ Array unidimensional com **índice**
- ▶ Como um dicionário (chave → valor)
- ▶ Como uma coluna de planilha
- ▶ Pode conter qualquer tipo de dados

Componentes:

- ▶ **Values:** Os dados propriamente ditos (array NumPy)
- ▶ **Index:** Rótulos para cada valor
- ▶ **Name:** Nome da Series (opcional)
- ▶ **Dtype:** Tipo de dados

💡 Nota Importante

Series = Array NumPy + Index

Criando Series: A partir de Lista

</> Python

```
1 # Series simples (indice automatico)
2 s = pd.Series([10, 20, 30, 40, 50])
3 print(s)
4 # 0      10
5 # 1      20
6 # 2      30
7 # 3      40
8 # 4      50
9 # dtype: int64
10
11 print(type(s)) # <class 'pandas.core.series.Series'>
12
```

Criando Series: Com Índice Customizado

</> Python

```
1 # Series com indice customizado
2 s = pd.Series([10, 20, 30, 40, 50],
3               index=['a', 'b', 'c', 'd', 'e'])
4 print(s)
5 # a      10
6 # b      20
7 # c      30
8 # d      40
9 # e      50
10 # dtype: int64
11
```

Criando Series: A partir de Dicionário

</> Python

```
1 # Dicionario -> Series
2 dados = {
3     'Ana': 25,
4     'Bruno': 30,
5     'Carlos': 35,
6     'Diana': 28
7 }
8 s = pd.Series(dados)
9 print(s)
10 # Ana      25
11 # Bruno    30
12 # Carlos   35
13 # Diana    28
14 # dtype: int64
15
```

Criando Series: Com Nome

</> Python

```
1 s = pd.Series([25, 30, 35, 28],
2               index=['Ana', 'Bruno', 'Carlos', 'Diana'],
3               name='Idade')
4 print(s)
5 # Ana      25
6 # Bruno    30
7 # Carlos   35
8 # Diana    28
9 # Name: Idade, dtype: int64
10
11 print(s.name) # 'Idade'
12
```


Parâmetros de Series

</> Python

```
1 pd.Series(  
2     data=None, # Lista, tupla, Array Numpy, etc.  
3     index=None, # Rótulo das linhas - IMUTÁVEL INDIVIDUALMENTE  
4     dtype=None, # Tipo de dados (um só)  
5     name=None, # Nome vira rótulo da coluna em um DF  
6     copy=None, # Copy-on-write (None: cópia preguiçosa)  
7     fastpath=False # Pula validação redundante (interno, instável)  
8 )
```

Atributos de Series

</> Python

```
1 s = pd.Series([10, 20, 30, 40],
2               index=['a', 'b', 'c', 'd'],
3               name='valores')
4
5 # Acessar valores (array NumPy)
6 print(s.values) # array([10, 20, 30, 40])
7 # Acessar indice
8 print(s.index)  # Index(['a', 'b', 'c', 'd'])
9 # Tipo de dados
10 print(s.dtype)  # int64 - Não é o tipo do índice!
11 # Tamanho
12 print(s.shape)  # (4,)
13 print(len(s))   # 4
```

Indexação em Series: Por Posição

</> Python

```
1 s = pd.Series([10, 20, 30, 40, 50],
2               index=['a', 'b', 'c', 'd', 'e'])
3
4 # Indexacao por posicao (como NumPy)
5 print(s[0])      # 10
6 print(s[-1])     # 50
7 print(s[1:3])    # Series com 'b' e 'c'
8
9 # Slicing
10 print(s[:3])     # Primeiros 3 elementos
11
```

Indexação em Series: Por Rótulo

Python

```
1 s = pd.Series([10, 20, 30, 40, 50],
2               index=['a', 'b', 'c', 'd', 'e'])
3 # Indexacao por rotulo (uma das principais vantagens do Pandas)
4 print(s['a'])          # 10
5 print(s['c'])          # 30
6 # Multiplos rotulos
7 print(s[['a', 'c', 'e']])
8 # a      10
9 # c      30
10 # e      50
```

Operações com Series

Python

```
1 s = pd.Series([10, 20, 30, 40, 50])
2
3 # Operacoes aritmeticas (vetorizadas)
4 print(s + 10)      # Adiciona 10 a cada elemento
5 print(s * 2)       # Multiplica por 2
6 print(s ** 2)      # Eleva ao quadrado
7
8 # Funcoes estatisticas
9 print(s.mean())    # 30.0
10 print(s.sum())     # 150
11 print(s.max())     # 50
12 print(s.std())     # ~15.81
13
```

Operações entre Series

</> Python

```
1 s1 = pd.Series([10, 20, 30], index=['a', 'b', 'c'])
2 s2 = pd.Series([1, 2, 3], index=['a', 'b', 'c'])
3 # Operacoes elemento a elemento
4 print(s1 + s2)
5 # a      11
6 # b      22
7 # c      33
8
9 # Alinhamento automatico por indice!
10 s3 = pd.Series([1, 2], index=['a', 'c'])
11 print(s1 + s3)
12 # a      11.0
13 # b       NaN
14 # c      33.0
```

DataFrame: Tabela 2D

O que é um DataFrame?

- ▶ Estrutura tabular bidimensional
- ▶ Como uma planilha do Excel
- ▶ Como uma tabela SQL
- ▶ Coleção de Series (cada coluna é uma Series)

Componentes:

- ▶ **Columns:** Nomes das colunas (Index)
- ▶ **Index:** Rótulos das linhas (Index)
- ▶ **Values:** Dados (array NumPy 2D)
- ▶ **Dtypes:** Tipo de cada coluna

💡 Nota Importante

DataFrame = Dicionário de Series (colunas)

Anatomia de um DataFrame

</> Python

```
1 import pandas as pd
2
3 df = pd.DataFrame({
4     'Nome': ['Ana', 'Bruno', 'Carlos'],
5     'Idade': [25, 30, 35],
6     'Cidade': ['SP', 'RJ', 'BH']
7 })
8
9 print(df)
10 #      Nome  Idade  Cidade
11 # 0     Ana    25     SP
12 # 1    Bruno    30     RJ
13 # 2   Carlos    35     BH
14
```


Criando DataFrames: A partir de Dicionário

</> Python

```
1 # Dicionario de listas
2 dados = {
3     'Nome': ['Ana', 'Bruno', 'Carlos', 'Diana'],
4     'Idade': [25, 30, 35, 28],
5     'Cidade': ['SP', 'RJ', 'BH', 'SP'],
6     'Salario': [3000, 4500, 5500, 3800]
7 }
8
9 df = pd.DataFrame(dados)
10 print(df)
11
```

Criando DataFrames: Resultado

</> Python

```
1 #      Nome  Idade  Cidade  Salario
2 # 0      Ana    25      SP    3000
3 # 1     Bruno   30      RJ    4500
4 # 2    Carlos   35      BH    5500
5 # 3     Diana   28      SP    3800
6
```

Criando DataFrames: A partir de Lista de Dicionários

Python

```
1 # Lista de dicionarios (cada dict = uma linha)
2 dados = [
3     {'Nome': 'Ana', 'Idade': 25, 'Cidade': 'SP'},
4     {'Nome': 'Bruno', 'Idade': 30, 'Cidade': 'RJ'},
5     {'Nome': 'Carlos', 'Idade': 35, 'Cidade': 'BH'}
6 ]
7
8 df = pd.DataFrame(dados)
9 print(df)
10 #      Nome  Idade  Cidade
11 # 0     Ana    25     SP
12 # 1    Bruno    30     RJ
13 # 2   Carlos    35     BH
14
```

Criando DataFrames: A partir de Array NumPy

Python

```
1 import numpy as np
2
3 # Array NumPy
4 arr = np.array([[1, 2, 3],
5                 [4, 5, 6],
6                 [7, 8, 9]])
7
8 df = pd.DataFrame(arr,
9                   columns=['A', 'B', 'C'],
10                  index=['x', 'y', 'z'])
11 print(df)
12 #      A  B  C
13 # x    1  2  3
14 # y    4  5  6
15 # z    7  8  9
```

Criando DataFrames: Com Índice Customizado

Python

```
1 dados = {  
2     'Nome': ['Ana', 'Bruno', 'Carlos'],  
3     'Idade': [25, 30, 35]  
4 }  
5  
6 df = pd.DataFrame(dados, index=['emp001', 'emp002', 'emp003'])  
7 print(df)  
8 #           Nome  Idade  
9 # emp001      Ana    25  
10 # emp002     Bruno    30  
11 # emp003     Carlos    35  
12
```

Leitura de Arquivos: CSV

Forma mais comum de criar DataFrames:

</> Python

```
1 # Ler CSV
2 df = pd.read_csv('titanic.csv') # extremamente flexível e poderoso
3
4 # Com opções
5 df = pd.read_csv('dados.csv',
6                  sep=';',          # Separador
7                  encoding='utf-8', # Codificação
8                  index_col=0,      # Coluna como índice
9                  na_values=['NA', '?']) # Valores faltantes
10
```

Leitura de Arquivos: Outros Formatos

Python

```
1 # Excel
2 df = pd.read_excel('dados.xlsx', sheet_name='Plan1')
3
4 # JSON
5 df = pd.read_json('dados.json')
6
7 # SQL
8 import sqlite3
9 conn = sqlite3.connect('database.db')
10 df = pd.read_sql('SELECT * FROM tabela', conn)
11
12 # Clipboard (copiar do Excel e colar)
13 df = pd.read_clipboard()
14
```

Atributos do DataFrame

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Shape (linhas, colunas)
3 print(df.shape) # (891, 12)
4
5 # Colunas
6 print(df.columns) # Index(['PassengerId', 'Survived', ...])
7
8 # Indice
9 print(df.index) # RangeIndex(start=0, stop=891, step=1)
10
11 # Tipos de dados
12 print(df.dtypes)
13
14 # Valores (array NumPy)
15 print(df.values) # array 2D
```


Parâmetros do DataFrame - set_index()

Python

```
1 DataFrame.set_index(  
2     keys,                # str, Index, Series, list, etc.  
3     drop=True,           # Remove do DF a coluna que virou index  
4     append=False,        # Adiciona novo index aos já existentes  
5     inplace=False,       # False: retorna view  
6     verify_integrity=False # Verifica índices duplicados  
7 )
```

Nota Importante

Para mais detalhes: Recursos Adicionais

DataFrame vs NumPy Array

NumPy Array:

</> Python

```
1 arr = np.array([[1, 2],
2                 [3, 4]])
3 # Indexacao por posicao
4 print(arr[0, 1]) # 2
5
6 # Sem rotulos
7 # Tipos homogeneos
8
```

DataFrame:

</> Python

```
1 df = pd.DataFrame([[1, 2],
2                     [3, 4]],
3                     columns=['A', 'B'])
4 # Indexacao por rotulo
5 print(df.loc[0, 'B']) # 2
6
7 # Com rotulos
8 # Tipos heterogeneos
9
```

Conversão: Pandas ↔ NumPy

Python

```
1 # DataFrame -> NumPy
2 df = pd.DataFrame([[1, 2], [3, 4]], columns=['A', 'B'])
3 arr = df.values # ou df.to_numpy()
4 print(type(arr)) # <class 'numpy.ndarray'>
5
6 # NumPy -> DataFrame
7 arr = np.array([[1, 2], [3, 4]])
8 df = pd.DataFrame(arr, columns=['A', 'B'])
9 print(type(df)) # <class 'pandas.core.frame.DataFrame'>
10
11 # Series -> NumPy
12 s = pd.Series([1, 2, 3])
13 arr = s.values
14
```

Relação: DataFrame e Series

</> Python

```
1 df = pd.DataFrame({
2     'Nome': ['Ana', 'Bruno', 'Carlos'],
3     'Idade': [25, 30, 35]
4 })
5
6 #DataFrame = Coleção de Series alinhadas pelo índice
7 # Cada coluna é uma Series
8 col = df['Nome']
9 print(type(col)) # <class 'pandas.core.series.Series'>
10
11 # Cada linha é uma Series
12 linha = df.loc[0]
13 print(type(linha)) # <class 'pandas.core.series.Series'>
```

Exemplo Prático: Carregar Titanic (arquivo local)

</> Python

```
1 import pandas as pd
2
3 # Carregar dataset
4 df = pd.read_csv('titanic.csv')
5
6 # Primeiras impressões
7 print(f"Linhas: {len(df)}")
8 print(f"Colunas: {len(df.columns)}")
9 print(f"Shape: {df.shape}")
10
11 # Tipos de dados presentes
12 print(df.dtypes.value_counts())
13
```

Exemplo Prático: Carregar Titanic (arquivo remoto)

</> Python

```
1 import pandas as pd
2
3 # Carregar dataset
4 url = "https://raw.githubusercontent.com/datasciencedojo/datasets/master/
      titanic.csv"
5 df = pd.read_csv(url)
6
7 # Primeiras impressoes
8 print(f"Linhas: {len(df)}")
9 print(f"Colunas: {len(df.columns)}")
10 print(f"Shape: {df.shape}")
11
12 # Tipos de dados presentes
13 print(df.dtypes.value_counts())
14
```

Bloco 2

Visualização e Inspeção de Dados

Explorando DataFrames

Primeiros passos com dados novos:

1. **Visualizar:** Primeiras/últimas linhas
2. **Dimensões:** Quantas linhas e colunas?
3. **Tipos:** Quais tipos de dados?
4. **Valores faltantes:** Existem NaN?
5. **Estatísticas:** Resumo numérico
6. **Distribuições:** Como são os valores?

Funções principais:

- ▶ `head()`, `tail()`, `sample()`
- ▶ `info()`, `describe()`
- ▶ `value_counts()`, `unique()`

head() e tail(): Visualização Rápida

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Primeiras 5 linhas (padrao)
4 print(df.head())
5
6 # Primeiras 10 linhas
7 print(df.head(10))
8
9 # Ultimas 5 linhas
10 print(df.tail())
11
12 # Ultimas 3 linhas
13 print(df.tail(3))
14
```

head() e tail(): Visualização Rápida

Nota Importante

Sempre comece com `head()` para entender a estrutura

sample(): Amostra Aleatória

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # 5 linhas aleatorias
4 print(df.sample(5))
5
6 # 10% dos dados
7 print(df.sample(frac=0.1))
8
9 # Com seed para reproducibilidade
10 print(df.sample(5, random_state=42))
11
12 # Amostragem com reposicao
13 print(df.sample(10, replace=True))
14
```

info(): Informações Gerais

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 df.info()
3 # <class 'pandas.core.frame.DataFrame'>
4 # RangeIndex: 891 entries, 0 to 890
5 # Data columns (total 12 columns):
6 #  #   Column      Non-Null Count  Dtype
7 #  ---  -
8 #  0   PassengerId  891 non-null    int64
9 #  1   Survived     891 non-null    int64
10 #  2   Pclass       891 non-null    int64
11 #  3   Name         891 non-null    object
12 #  4   Sex          891 non-null    object
13 #  5   Age         714 non-null    float64
14 #  ...
```

info(): O que mostra

Informações fornecidas:

- ▶ Tipo da classe (DataFrame)
- ▶ Tipo e tamanho do índice
- ▶ Número total de colunas
- ▶ Para cada coluna:
 - ▶ Nome
 - ▶ Número de valores não-nulos
 - ▶ Tipo de dados (dtype)
- ▶ Uso de memória

💡 Nota Importante

`info()` é essencial para identificar valores faltantes rapidamente

describe(): Estatísticas Descritivas

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Estatísticas para colunas numericas
4 print(df.describe())
5 #
6 # count      PassengerId      Survived      Pclass      Age
7 # mean        446.000000      0.383838      2.308642     29.699118
8 # std         257.353842      0.486592      0.836071     14.526497
9 # min          1.000000      0.000000      1.000000      0.420000
10 # 25%         223.500000      0.000000      2.000000     20.125000
11 # 50%         446.000000      0.000000      3.000000     28.000000
12 # 75%         668.500000      1.000000      3.000000     38.000000
13 # max         891.000000      1.000000      3.000000     80.000000
14
```

describe(): Incluir Todas as Colunas

Python

```
1 # Incluir colunas categoricas
2 print(df.describe(include='all'))
3
4 # Apenas colunas categoricas
5 print(df.describe(include='object'))
6
7 # Apenas colunas numericas (padrao)
8 print(df.describe(include='number'))
9
10 # Estatisticas especificas
11 print(df[['Age', 'Fare']].describe())
12
```

describe(): Interpretação

Para colunas numéricas:

- ▶ **count:** Número de valores não-nulos
- ▶ **mean:** Média
- ▶ **std:** Desvio padrão
- ▶ **min/max:** Valores mínimo e máximo
- ▶ **25%/50%/75%:** Percentis (quartis)

Para colunas categóricas:

- ▶ **count:** Número de valores não-nulos
- ▶ **unique:** Número de valores únicos
- ▶ **top:** Valor mais frequente
- ▶ **freq:** Frequência do valor mais frequente

value_counts(): Distribuição de Valores

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Contar valores unicos
3 print(df['Survived'].value_counts())
4 # 0      549
5 # 1      342
6 # Name: Survived, dtype: int64
7
8 # Com proporcoes (percentuais)
9 print(df['Survived'].value_counts(normalize=True))
10 # 0      0.616162
11 # 1      0.383838
12
13 # Incluir valores faltantes
14 print(df['Age'].value_counts(dropna=False))
```

value_counts(): Aplicações

</> Python

```
1 # Distribuicao por classe
2 print(df['Pclass'].value_counts().sort_index())
3 # 1      216
4 # 2      184
5 # 3      491
6
7 # Distribuicao por sexo
8 print(df['Sex'].value_counts())
9 # male      577
10 # female    314
11
12 # Top 5 nomes mais comuns (se houver repetidos)
13 print(df['Name'].value_counts().head(5))
14
```

unique() e nunique()

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Valores unicos (array)
3 print(df['Pclass'].unique())
4 # array([3, 1, 2])
5
6 # Numero de valores unicos
7 print(df['Pclass'].nunique())
8 # 3
9
10 # Quantos nomes unicos?
11 print(df['Name'].nunique())
12 # 891 (todos diferentes)
13
14 # Quantas idades unicas?
15 print(df['Age'].nunique())
```

Atributos: shape, columns, index

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Shape (linhas, colunas)
4 print(df.shape) # (891, 12)
5 print(f"Linhas: {df.shape[0]}")
6 print(f"Colunas: {df.shape[1]}")
7
8 # Nomes das colunas
9 print(df.columns)
10 # Index(['PassengerId', 'Survived', 'Pclass', ...])
11
12 # Indice (linhas)
13 print(df.index)
14 # RangeIndex(start=0, stop=891, step=1)
15
```

dtypes: Tipos de Dados

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Tipos de cada coluna
3 print(df.dtypes)
4 # PassengerId      int64
5 # Survived         int64
6 # Pclass           int64
7 # Name             object
8 # Sex              object
9 # Age              float64
10 # ...
11 # Resumo dos tipos
12 print(df.dtypes.value_counts())
13 # int64           5
14 # float64         2
15 # object          5
```

Tipos de Dados no Pandas

Principais dtypes:

Dtype	Descrição
int64	Inteiros
float64	Ponto flutuante
object	Strings, tipos mistos
bool	Booleanos (True/False)
datetime64	Datas e horários
timedelta64	Diferenças de tempo
category	Categóricas (economia de memória)

💡 Nota Importante

object geralmente significa strings, mas pode conter qualquer tipo Python

Verificar Valores Faltantes

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Valores faltantes por coluna
3 print(df.isnull().sum())
4 # PassengerId      0
5 # Survived         0
6 # Pclass           0
7 # Name             0
8 # Sex              0
9 # Age             177
10 # Cabin           687
11 # Embarked        2
12 # ...
13
14 # Percentual de valores faltantes
15 print((df.isnull().sum() / len(df)) * 100)
```

Verificar Valores Faltantes: Total

</> Python

```
1 # Total de valores faltantes no DataFrame
2 print(df.isnull().sum().sum())
3
4 # Linhas com pelo menos um valor faltante
5 print(df.isnull().any(axis=1).sum())
6
7 # Linhas sem valores faltantes
8 print(df.notnull().all(axis=1).sum())
9
10 # Apenas linhas completas (sem NaN)
11 df_completo = df.dropna()
12 print(f"Linhas completas: {len(df_completo)}")
13
```


Renomear Colunas: rename()

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Renomear colunas especificas
3 df_renamed = df.rename(columns={
4     'PassengerId': 'ID',
5     'Pclass': 'Classe',
6     'Name': 'Nome'
7 })
8
9 # Renomear in-place (modifica o original)
10 df.rename(columns={'Sex': 'Sexo'}, inplace=True)
11
12 # Aplicar funcao a todos os nomes
13 df.columns = df.columns.str.lower() # Minusculas
14 df.columns = df.columns.str.replace(' ', '_')
```

Renomear Índice: set_index()

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Definir coluna como índice
3 df_indexed = df.set_index('PassengerId')
4 print(df_indexed.head())
5 #           Survived  Pclass  Name  ...
6 # PassengerId
7 # 1              0      3  ...  ...
8 # 2              1      1  ...  ...
9
10 # Resetar índice
11 df_reset = df_indexed.reset_index()
12
13 # set_index in-place
14 df.set_index('Name', inplace=True)
```

Ordenar Colunas

Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Reordenar colunas manualmente
4 colunas_ordem = ['Name', 'Age', 'Sex', 'Survived']
5 df_reord = df[colunas_ordem]
6
7 # Ordenar alfabeticamente
8 df_sorted = df[sorted(df.columns)]
9
10 # Mover coluna para o inicio
11 cols = df.columns.tolist()
12 cols = ['Survived'] + [c for c in cols if c != 'Survived']
13 df = df[cols]
14
```

Exemplo Prático: Exploração Inicial

</> Python

```
1 import pandas as pd
2
3 # Carregar dados
4 df = pd.read_csv('titanic.csv')
5
6 # Pipeline de exploracao inicial
7 print("=== DIMENSOES ===")
8 print(f"Shape: {df.shape}")
9
10 print("\n=== PRIMEIRAS LINHAS ===")
11 print(df.head(3))
12
13 print("\n=== INFORMACOES ===")
14 df.info()
15
```

Exemplo Prático: Exploração (cont.)

Python

```
1 print("\n=== ESTATISTICAS ===")
2 print(df.describe())
3
4 print("\n=== VALORES FALTANTES ===")
5 print(df.isnull().sum())
6
7 print("\n=== DISTRIBUICAO DE SOBREVIVENTES ===")
8 print(df['Survived'].value_counts())
9
10 print("\n=== DISTRIBUICAO POR CLASSE ===")
11 print(df['Pclass'].value_counts().sort_index())
12
```

Bloco 3

Seleção e Indexação

Formas de Selecionar Dados

Pandas oferece múltiplas formas:

- ▶ **Colchetes []**: Seleção de colunas
- ▶ **.loc[]**: Indexação por **rótulo**
- ▶ **.iloc[]**: Indexação por **posição**
- ▶ **.at[]**: Acesso rápido a valor único (rótulo)
- ▶ **.iat[]**: Acesso rápido a valor único (posição)

Quando usar cada um?

para colunas

- ▶ **.loc** para seleção por nomes
- ▶ **.iloc** para seleção por posição numérica
- ▶ **.at/.iat** para acesso a elemento único (mais rápido)

Seleção de Colunas: Uma Coluna

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Selecionar uma coluna (retorna Series)
4 nomes = df['Name']
5 print(type(nomes)) # <class 'pandas.core.series.Series'>
6
7 # Alternativa (notacao de atributo)
8 # Funciona se nome nao tem espacos/caracteres especiais
9 idades = df.Age # Mesmo que df['Age']
10
11 # Ver alguns valores
12 print(nomes.head())
13
```


Seleção de Colunas: Múltiplas Colunas

Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Selecionar multiplas colunas (retorna DataFrame)
4 subset = df[['Name', 'Age', 'Sex']]
5 print(type(subset)) # <class 'pandas...DataFrame'>
6 print(subset.head())
7
8 # Ordem importa
9 subset2 = df[['Sex', 'Name', 'Age']]
10
11 # Todas exceto algumas
12 colunas_manter = [c for c in df.columns
13                   if c not in ['PassengerId', 'Ticket']]
14 df_filtrado = df[colunas_manter]
15
```

Diferença: [] vs [[]]

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Uma coluna: retorna Series
3 series = df['Age']
4 print(type(series)) # Series
5 print(series.shape) # (891,)
6
7 # Lista com uma coluna: retorna DataFrame
8 dataframe = df[['Age']]
9 print(type(dataframe)) # DataFrame
10 print(dataframe.shape) # (891, 1)
11
12 # Diferença visual
13 print(series.head(2))      # 0      22.0
14                          # 1      38.0
15 print(dataframe.head(2)) # DataFrame formatado
```

loc vs iloc: Conceito

Diferença fundamental:

.loc[]:

- ▶ Indexação por **rótulo**
- ▶ Usa nomes de índice/colunas
- ▶ Inclusivo nos extremos
- ▶ `df.loc[0:5]` inclui linha 5

.iloc[]:

- ▶ Indexação por **posição**
- ▶ Usa números inteiros
- ▶ Exclusivo no extremo direito
- ▶ `df.iloc[0:5]` exclui linha 5

💡 Nota Importante

loc = location (rótulo), **iloc** = integer location (posição)

loc[]: Seleção por Rótulo

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Uma linha por índice
4 print(df.loc[0]) # Series com dados da linha 0
5
6 # Múltiplas linhas
7 print(df.loc[0:5]) # Linhas 0, 1, 2, 3, 4, 5 (inclusivo!)
8
9 # Uma célula específica
10 print(df.loc[0, 'Name'])
11
12 # Subconjunto (linhas e colunas)
13 subset = df.loc[0:5, ['Name', 'Age', 'Sex']]
14 print(subset)
15
```

loc[]: Slicing de Colunas

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Todas as linhas, algumas colunas
4 subset = df.loc[:, ['Name', 'Age', 'Sex']]
5
6 # Slice de colunas (por nome)
7 subset = df.loc[:, 'Name':'Age'] # Inclusivo!
8
9 # Primeiras 10 linhas, colunas específicas
10 subset = df.loc[:9, ['Survived', 'Pclass', 'Name']]
11
12 # Uma linha, todas as colunas
13 linha = df.loc[5, :]
14
```

iloc[]: Seleção por Posição

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Uma linha por posicao
4 print(df.iloc[0]) # Primeira linha
5
6 # Multiplas linhas
7 print(df.iloc[0:5]) # Linhas 0, 1, 2, 3, 4 (exclusivo!)
8
9 # Uma celula especifica (linha 0, coluna 3)
10 print(df.iloc[0, 3])
11
12 # Subconjunto (primeiras 5 linhas, primeiras 3 colunas)
13 subset = df.iloc[0:5, 0:3]
14 print(subset)
15
```

iloc[:]: Indexação Negativa

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Última linha
3 print(df.iloc[-1])
4
5 # Últimas 5 linhas
6 print(df.iloc[-5:])
7
8 # Penúltima linha, últimas 3 colunas
9 print(df.iloc[-2, -3:])
10
11 # Todas as linhas exceto a primeira
12 print(df.iloc[1:])
13
14 # Primeiras 100 linhas, exceto primeiras 2 colunas
15 subset = df.iloc[:100, 2:]
```

iloc[]: Listas de Posições

Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Linhas especificas
4 linhas = df.iloc[[0, 10, 20, 30]]
5
6 # Colunas especificas (posicoes 1, 3, 5)
7 colunas = df.iloc[:, [1, 3, 5]]
8
9 # Combinacao
10 subset = df.iloc[[0, 5, 10], [1, 3, 5]]
11
12 # Linhas pares
13 pares = df.iloc[:, :2] # step=2
14
```


loc vs iloc: Comparação Direta

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # loc: por rotulo (inclusivo)
4 a = df.loc[0:5, 'Name']
5 print(len(a)) # 6 elementos (0, 1, 2, 3, 4, 5)
6
7 # iloc: por posicao (exclusivo)
8 b = df.iloc[0:5, 3] # coluna na posicao 3
9 print(len(b)) # 5 elementos (0, 1, 2, 3, 4)
10
11 # loc com indice customizado
12 df_custom = df.set_index('PassengerId')
13 print(df_custom.loc[1]) # Passageiro ID 1
14 print(df_custom.iloc[0]) # Primeira linha (pode ser != 1)
15
```

at[] e iat[]: Acesso Rápido

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # at: acesso por rotulo (mais rapido que loc)
4 valor = df.at[0, 'Name']
5 print(valor) # Nome do passageiro 0
6
7 # iat: acesso por posicao (mais rapido que iloc)
8 valor = df.iat[0, 3]
9 print(valor) # Valor na linha 0, coluna 3
10
11 # Modificar valor
12 df.at[0, 'Age'] = 25
13 df.iat[0, 5] = 30 # Se Age estiver na posicao 5
14
```

Quando Usar Cada Método

Método	Uso	Retorna
<code>df['col']</code>	Uma coluna	Series
<code>df[['c1', 'c2']]</code>	Múltiplas colunas	DataFrame
<code>df.loc[]</code>	Rótulos (nomes)	Series/DataFrame
<code>df.iloc[]</code>	Posições (int)	Series/DataFrame
<code>df.at[]</code>	Um valor (rótulo)	Escalar
<code>df.iat[]</code>	Um valor (posição)	Escalar

Regra geral:

- ▶ Colunas: `df['col']` ou `df[['col1', 'col2']]`
- ▶ Linhas e colunas: `.loc` ou `.iloc`
- ▶ Um valor: `.at` ou `.iat`

Modificando Valores

Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Modificar uma coluna inteira
4 df['Age'] = df['Age'] + 1
5
6 # Modificar uma célula específica
7 df.loc[0, 'Age'] = 25
8 df.iloc[0, 5] = 30
9
10 # Modificar múltiplas células
11 df.loc[0:5, 'Age'] = 30
12
13 # Criar nova coluna
14 df['AgeGroup'] = 'Adult'
15
```

Localizando Valores Aninhados

Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Potencial problema (chained indexing)
4 # df['Age'][0] = 25 # Pode dar warning!
5
6 # Forma correta
7 df.loc[0, 'Age'] = 25
8
9 # Criar copia explicita se necessario
10 subset = df[['Name', 'Age']].copy()
11 subset.loc[0, 'Age'] = 25 # OK, é uma copia
12
```

! Chained Indexing

⚠ Atenção

Evite **chained indexing**: `df['col'][row]`. Use `df.loc[row, 'col']`

Exemplo Prático: Extrair Subset

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Extrair informacoes basicas dos primeiros 100
4 subset = df.loc[:99, ['Name', 'Age', 'Sex', 'Survived']]
5 print(subset.head())
6
7 # Extrair ultimos 50 passageiros, todas as colunas
8 ultimos = df.iloc[-50:]
9
10 # Extrair colunas especificas por posicao
11 # Colunas 1, 2, 3, 5 (Survived, Pclass, Name, Sex)
12 especificas = df.iloc[:, [1, 2, 3, 5]]
13
```

Bloco 4

Filtragem e Ordenação

Filtragem de Dados

Por que filtrar?

- ▶ Análise de subconjuntos específicos
- ▶ Remover outliers ou dados inválidos
- ▶ Focar em categorias de interesse
- ▶ Preparar dados para modelagem

Métodos principais:

- ▶ **Boolean indexing:** Máscaras booleanas
- ▶ **Operadores lógicos:** `&`, `|`, `~`
- ▶ **query():** Sintaxe SQL-like
- ▶ **isin():** Verificar pertencimento

Boolean Indexing: Conceito

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Criar mascara booleana
4 mascara = df['Age'] > 30
5 print(type(mascara)) # Series de booleanos
6 print(mascara.head())
7 # 0      False
8 # 1       True
9 # 2      False
10 # 3       True
11 # 4      False
12 # Aplicar mascara
13 maiores_30 = df[mascara]
14 print(len(maiores_30))
```

Boolean Indexing: Forma Compacta

Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Forma compacta (mais comum)
4 maiores_30 = df[df['Age'] > 30]
5
6 # Outras comparacoes
7 primeira_classe = df[df['Pclass'] == 1]
8 mulheres = df[df['Sex'] == 'female']
9 sobreviventes = df[df['Survived'] == 1]
10
11 # Combinacao comum
12 mulheres_sobreviventes = df[(df['Sex'] == 'female') &
13                               (df['Survived'] == 1)]
14
```

Operadores de Comparação

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Maior que
3 df[df['Age'] > 30]
4
5 # Menor ou igual
6 df[df['Fare'] <= 20]
7
8 # Igual
9 df[df['Pclass'] == 1]
10
11 # Diferente
12 df[df['Sex'] != 'male']
13
14 # Entre valores
15 df[(df['Age'] >= 18) & (df['Age'] <= 35)]
```

Operadores Lógicos: AND (&)

Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # AND: ambas condicoes devem ser True
4 mulheres_1classe = df[(df['Sex'] == 'female') &
5                        (df['Pclass'] == 1)]
6 print(len(mulheres_1classe))
7
8 # Multiplas condicoes
9 filtro = df[(df['Age'] > 30) &
10             (df['Pclass'] == 1) &
11             (df['Survived'] == 1)]
12
```

Operadores Lógicos: AND (&) (cont.)

</> Python

```
1 # Equivalente (mais legível para muitas condições)
2 cond1 = df['Age'] > 30
3 cond2 = df['Pclass'] == 1
4 cond3 = df['Survived'] == 1
5 filtro = df[cond1 & cond2 & cond3]
6
```

Operadores Lógicos: OR (|)

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # OR: pelo menos uma condicao deve ser True
4 primeira_ou_segunda = df[(df['Pclass'] == 1) |
5                           (df['Pclass'] == 2)]
6 print(len(primeira_ou_segunda))
7
8 # Crianças ou idosos
9 extremos = df[(df['Age'] < 12) | (df['Age'] > 60)]
10
11 # Combinacao AND e OR
12 filtro = df[((df['Sex'] == 'female') | (df['Age'] < 18)) &
13             (df['Survived'] == 1)]
14
```

Operadores Lógicos: OR (|)

Atenção

Use & e |, não and e or! Sempre use parênteses!

Operadores Lógicos: NOT (~)

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # NOT: inverte a condicao
3 nao_sobreviventes = df[~(df['Survived'] == 1)]
4 # Equivalente a:
5 nao_sobreviventes = df[df['Survived'] == 0]
6
7 # Nao-criancas (idade >= 18)
8 adultos = df[~(df['Age'] < 18)]
9
10 # Excluir primeira classe
11 nao_primeira = df[~(df['Pclass'] == 1)]
12
13 # Combinacao complexa
14 filtro = df[~((df['Age'] < 18) | (df['Age'] > 60))]
15 # Pessoas entre 18 e 60 anos
```

isin(): Pertencimento à Lista

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Verificar se valor esta em lista
3 portos = ['S', 'C']
4 embarque_filtrado = df[df['Embarked'].isin(portos)]
5
6 # Classes 1 ou 2
7 classes_altas = df[df['Pclass'].isin([1, 2])]
8
9 # NOT isin (nao esta na lista)
10 nao_Southampton = df[~df['Embarked'].isin(['S'])]
11
12 # Nomes especificos
13 nomes = ['Smith', 'Johnson', 'Williams']
14 pessoas = df[df['Name'].str.contains('|'.join(nomes))]
```

query(): Sintaxe SQL-like

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Sintaxe query (string)
3 resultado = df.query('Age > 30')
4
5 # Multiplas condicoes
6 resultado = df.query('Age > 30 and Pclass == 1')
7 # OR
8 resultado = df.query('Pclass == 1 or Pclass == 2')
9 # Variavel externa
10 idade_minima = 30
11 resultado = df.query('Age > @idade_minima')
12
13 # Nomes com espacos (usar backticks)
14 # df.query('`Passenger Class` == 1')
15
```

query(): Vantagens

Vantagens do query():

- ▶ Sintaxe mais legível para queries complexas
- ▶ Não precisa repetir `df['coluna']`
- ▶ Usa `and/or` em vez de `&/|`
- ▶ Não precisa de parênteses extras

query(): Vantagens

</> Python

```
1 # Comparacao: boolean indexing vs query
2 # Boolean indexing (verboso)
3 filtro1 = df[(df['Age'] > 30) & (df['Pclass'] == 1)]
4
5 # Query (conciso)
6 filtro2 = df.query('Age > 30 and Pclass == 1')
7
8 # Resultados identicos
9 print(filtro1.equals(filtro2)) # True
10
```

Filtragem com Strings

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Contem substring
3 mrs = df[df['Name'].str.contains('Mrs')]
4 # Comeca com
5 comeca_a = df[df['Name'].str.startswith('A')]
6 # Termina com
7 termina_son = df[df['Name'].str.endswith('son')]
8 # Case insensitive
9 braund = df[df['Name'].str.contains('braund',
10                                     case=False)]
11 # Regex
12 import re
13 tem_numero = df[df['Name'].str.contains(r'\d')]
```

Filtragem de Valores Faltantes

Python

```
1 df = pd.read_csv('titanic.csv')
2 # Linhas com idade faltante
3 sem_idade = df[df['Age'].isnull()]
4 print(len(sem_idade))
5
6 # Linhas com idade presente
7 com_idade = df[df['Age'].notnull()]
8 # ou: df[~df['Age'].isnull()]
9
10 # Linhas sem nenhum valor faltante
11 completas = df[df.notnull().all(axis=1)]
12
13 # Linhas com pelo menos um valor faltante
14 com_nan = df[df.isnull().any(axis=1)]
```

between(): Valores em Range

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Idade entre 20 e 40 (inclusivo)
4 faixa = df[df['Age'].between(20, 40)]
5
6 # Equivalente a:
7 faixa = df[(df['Age'] >= 20) & (df['Age'] <= 40)]
8
9 # Excluir extremos (exclusive)
10 faixa = df[df['Age'].between(20, 40, inclusive='neither')]
11
12 # Tarifa entre valores
13 tarifas = df[df['Fare'].between(10, 50)]
14
```


Ordenação de Dados

Por que ordenar?

- ▶ Identificar valores extremos
- ▶ Preparar dados para análise sequencial
- ▶ Facilitar visualização
- ▶ Ranking de elementos

Métodos principais:

- ▶ `sort_values()`: Ordenar por valores de coluna(s)
- ▶ `sort_index()`: Ordenar pelo índice
- ▶ `nlargest()/nsmallest()`: Top N valores

sort_values(): Ordenar por Coluna

Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Ordenar por idade (crescente)
4 df_ordenado = df.sort_values('Age')
5 print(df_ordenado.head())
6
7 # Ordenar por idade (decrescente)
8 df_ordenado = df.sort_values('Age', ascending=False)
9
10 # Ordenar in-place (modifica original)
11 df.sort_values('Age', inplace=True)
12
13 # Resetar indice apos ordenar
14 df_ordenado = df.sort_values('Age').reset_index(drop=True)
```

sort_values(): Múltiplas Colunas

Python

```
1 df = pd.read_csv('titanic.csv')
2 # Ordenar por classe, depois por tarifa
3 df_ordenado = df.sort_values(['Pclass', 'Fare'])
4
5 # Direcoes diferentes para cada coluna
6 df_ordenado = df.sort_values(
7     ['Pclass', 'Age'],
8     ascending=[True, False] # Classe crescente, Idade decrescente
9 )
10 # Exemplo: classe crescente, dentro de cada classe
11 # ordenar por tarifa decrescente
12 df_ordenado = df.sort_values(
13     ['Pclass', 'Fare'],
14     ascending=[True, False]
15 )
```

sort_values(): Tratando NaN

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Por padrao, NaN vai para o final
4 df_ordenado = df.sort_values('Age')
5
6 # Colocar NaN no inicio
7 df_ordenado = df.sort_values('Age', na_position='first')
8
9 # Remover NaN antes de ordenar
10 df_sem_nan = df.dropna(subset=['Age'])
11 df_ordenado = df_sem_nan.sort_values('Age')
12
13 # Ordenar apenas valores nao-nulos
14 df_ordenado = df[df['Age'].notnull()].sort_values('Age')
15
```

sort_index(): Ordenar por Índice

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Depois de filtros, indice pode estar desordenado
4 filtrado = df[df['Age'] > 30]
5 print(filtrado.index[:5]) # [1, 3, 6, 11, 15, ...]
6
7 # Ordenar pelo indice
8 df_ordenado = filtrado.sort_index()
9 print(df_ordenado.index[:5]) # [1, 3, 6, 11, 15, ...]
10
11 # Decrescente
12 df_ordenado = filtrado.sort_index(ascending=False)
13
14 # Para indice customizado (strings)
15 df_custom = df.set_index('Name').sort_index()
16
```

nlargest() e nsmallest(): Top N

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Top 10 tarifas mais caras
4 top_tarifas = df.nlargest(10, 'Fare')
5 print(top_tarifas[['Name', 'Fare']])
6
7 # 10 passageiros mais jovens
8 mais_jovens = df.nsmallest(10, 'Age')
9
10 # Top 5 por multiplas colunas
11 # (primeiro por Fare, depois por Age)
12 top_5 = df.nlargest(5, ['Fare', 'Age'])
13
```

nlargest() e nsmallest(): Top N

💡 Nota Importante

Mais eficiente que `sort_values().head(n)` para Top N

Exemplo Prático: Análise de Sobreviventes

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Mulheres que sobreviveram, primeira classe
4 filtro = df[(df['Sex'] == 'female') &
5             (df['Survived'] == 1) &
6             (df['Pclass'] == 1)]
7
8 # Ordenar por idade
9 filtro_ordenado = filtro.sort_values('Age')
10
11 # Ver informacoes basicas
12 resultado = filtro_ordenado[['Name', 'Age', 'Fare']]
13 print(f"Total: {len(resultado)}")
14 print(resultado.head(10))
15
```


Exemplo Prático: Top Passageiros

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # 10 passageiros que pagaram tarifas mais altas
4 top_10_tarifas = df.nlargest(10, 'Fare')
5
6 # Dessas, quantas sobreviveram?
7 sobreviventes = top_10_tarifas['Survived'].sum()
8 taxa = (sobreviventes / 10) * 100
9
10 print(f"Sobreviventes entre top 10 tarifas: {sobreviventes}")
11 print(f"Taxa de sobrevivencia: {taxa}%")
12
13 # Ver detalhes
14 print(top_10_tarifas[['Name', 'Fare', 'Survived']])
15
```

Exemplo Prático: Crianças por Classe

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Definir crianças (idade < 18)
3 crianas = df[df['Age'] < 18]
4
5 print(f"Total de crianças: {len(crianças)}")
6
7 # Distribuicao por classe
8 print("\nCrianças por classe:")
9 print(crianças['Pclass'].value_counts().sort_index())
10
11 # Taxa de sobrevivencia de crianças
12 taxa = crianças['Survived'].mean() * 100
13 print(f"\nTaxa sobrevivencia crianças: {taxa:.1f}%")
```

Exemplo Prático: Análise Complexa

Python

```
1 df = pd.read_csv('titanic.csv')
2 # Adultos (>= 18), primeira ou segunda classe
3 # que sobreviveram, ordenados por idade
4 resultado = df.query(
5     'Age >= 18 and (Pclass == 1 or Pclass == 2) '
6     'and Survived == 1'
7 ).sort_values('Age')
8 # Estatísticas
9 print(f"Total: {len(resultado)}")
10 print(f"Idade media: {resultado['Age'].mean():.1f}")
11 print(f"Tarifa media: {resultado['Fare'].mean():.2f}")
12 # Top 5 mais jovens
13 print("\n5 mais jovens:")
14 print(resultado.head(5)[['Name', 'Age', 'Pclass']])
```

Pipeline de Análise Completo

</> Python

```
1 import pandas as pd
2 # Carregar
3 df = pd.read_csv('titanic.csv')
4 # Filtrar
5 adultos_1classe = df[
6     (df['Age'] >= 18) &
7     (df['Pclass'] == 1)
8 ]
9 # Ordenar
10 ordenado = adultos_1classe.sort_values(
11     ['Survived', 'Age'],
12     ascending=[False, True]
13 )
14 # Selecionar colunas
15 resultado = ordenado[['Name', 'Age', 'Fare', 'Survived']]
```

Pipeline: Forma Encadeada (Method Chaining)

</> Python

```
1 import pandas as pd
2
3 # Pipeline completo em uma linha
4 resultado = (
5     pd.read_csv('titanic.csv')
6     .query('Age >= 18 and Pclass == 1')
7     .sort_values(['Survived', 'Age'],
8                 ascending=[False, True])
9     [['Name', 'Age', 'Fare', 'Survived']]
10 )
11
12 print(resultado.head())
13
```

Pipeline: Forma Encadeada (Method Chaining)

Nota Importante

Method chaining torna o código mais legível e conciso

Estatísticas de Subgrupos

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Estatísticas por grupo
4 print("=== PRIMEIRA CLASSE ===")
5 primeira = df[df['Pclass'] == 1]
6 print(f"Passageiros: {len(primeira)}")
7 print(f"Taxa sobrevivencia: {primeira['Survived'].mean():.2%}")
8 print(f"Idade media: {primeira['Age'].mean():.1f}")
9
10 print("\n=== TERCEIRA CLASSE ===")
11 terceira = df[df['Pclass'] == 3]
12 print(f"Passageiros: {len(terceira)}")
13 print(f"Taxa sobrevivencia: {terceira['Survived'].mean():.2%}")
14 print(f"Idade media: {terceira['Age'].mean():.1f}")
15
```

Criando Categorias (Binning Manual)

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Criar coluna de faixa etaria
3 df['FaixaEtaria'] = 'Adulto'
4 df.loc[df['Age'] < 18, 'FaixaEtaria'] = 'Crianca'
5 df.loc[df['Age'] >= 60, 'FaixaEtaria'] = 'Idoso'
6
7 # Ver distribuicao
8 print(df['FaixaEtaria'].value_counts())
9
10 # Analise por faixa
11 for faixa in ['Crianca', 'Adulto', 'Idoso']:
12     subset = df[df['FaixaEtaria'] == faixa]
13     taxa = subset['Survived'].mean()
14     print(f"{faixa}: {taxa:.2%}")
```


Salvar Resultados Filtrados

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Filtrar dados de interesse
3 sobreviventes = df[df['Survived'] == 1]
4 # Salvar em novo CSV
5 sobreviventes.to_csv('sobreviventes.csv', index=False)
6 # Salvar apenas colunas especificas
7 colunas = ['Name', 'Age', 'Sex', 'Pclass']
8 sobreviventes[colunas].to_csv(
9     'sobreviventes_basico.csv',
10     index=False
11 )
12 # Salvar em Excel
13 sobreviventes.to_excel('sobreviventes.xlsx',
14                         index=False,
15                         sheet_name='Sobreviventes')
```

Live Coding

Vamos praticar juntos:

1. Carregar dataset Titanic
2. Exploração inicial (head, info, describe)
3. Seleção de colunas (loc, iloc)
4. Filtragem com boolean indexing
5. Ordenação por múltiplas colunas
6. Análise de subgrupos

Bloco 1 - Estruturas Pandas:

- ▶ Series: array 1D com índice
- ▶ DataFrame: tabela 2D com rótulos
- ▶ Criação: dicionários, listas, arrays, CSV
- ▶ Relação com NumPy
- ▶ Atributos: shape, columns, index, dtypes

Bloco 2 - Visualização e Inspeção:

- ▶ head(), tail(), sample()
- ▶ info(), describe(), value_counts()
- ▶ Verificação de valores faltantes
- ▶ Renomeação: rename(), set_index()

Revisão da Aula (cont.)

Bloco 3 - Seleção e Indexação:

- ▶ Seleção de colunas: [] e [[]]
- ▶ loc[]: indexação por rótulo (inclusivo)
- ▶ iloc[]: indexação por posição (exclusivo)
- ▶ at[]/iat[]: acesso rápido a valor único
- ▶ SettingWithCopyWarning: evitar chained indexing

Bloco 4 - Filtragem e Ordenação:

- ▶ Boolean indexing: máscaras booleanas
- ▶ Operadores lógicos: &, |, ~
- ▶ query(): sintaxe SQL-like
- ▶ sort_values(), sort_index()
- ▶ nlargest(), nsmallest()

Conceitos-Chave para Lembrar

1. **Series = Array + Index**, DataFrame = Coleção de Series
2. **loc usa rótulos** (inclusivo), **iloc** usa posições (exclusivo)
3. **head()/info()/describe()** são essenciais para exploração inicial
4. **Boolean indexing**: `df[df['col'] > valor]`
5. **Use &, |, ~**, não `and`, `or`, `not`
6. **Sempre use parênteses** em condições múltiplas
7. **query()** é mais legível para filtros complexos
8. **sort_values()** para ordenar, **nlargest()** para Top N
9. **Method chaining** torna código mais limpo

Pandas vs NumPy: Quando Usar

Use NumPy quando:

- ▶ Performance crítica
- ▶ Arrays numéricos puros
- ▶ Álgebra linear
- ▶ Processamento de imagens
- ▶ Algoritmos matemáticos

Use Pandas quando:

- ▶ Dados tabulares (CSV, Excel)
- ▶ Colunas com tipos diferentes
- ▶ Análise exploratória
- ▶ Limpeza de dados
- ▶ Agregações e groupby
- ▶ Séries temporais

💡 Nota Importante

Pandas é construído sobre NumPy - use o melhor de cada um!

Ao trabalhar com Pandas:

1. **Sempre explore primeiro:** `head()`, `info()`, `describe()`
2. **Verifique valores faltantes:** `isnull().sum()`
3. **Use nomes descritivos:** `df_filtrado`, `sobreviventes`
4. **Comente filtros complexos:** o que você está filtrando?
5. **Use method chaining:** código mais limpo
6. **Evite loops:** Pandas já é vetorizado
7. **Use `copy()` quando necessário:** evitar modificações acidentais
8. **Documente transformações:** rastreabilidade
9. **Salve checkpoints:** `to_csv()` para dados intermediários

Erros Comuns e Como Evitar

1. **SettingWithCopyWarning**

- ▶ Evite: `df['col'][row] = value`
- ▶ Use: `df.loc[row, 'col'] = value`

2. **KeyError: coluna não existe**

- ▶ Verifique: `df.columns`
- ▶ Cuidado com espaços e case-sensitive

3. **Usar and/or em vez de &/|**

- ▶ Sempre use: `&`, `|`, `~`
- ▶ Sempre use parênteses: `(cond1) & (cond2)`

4. **Confundir loc com iloc**

- ▶ `loc` = labels (nomes)
- ▶ `iloc` = integer location (posições)

Próxima Aula: Manipulação de Dados

Aula 06 - O que vem:

- ▶ Transformações: `apply()`, `map()`, `applymap()`
- ▶ Operações com strings
- ▶ `GroupBy` e agregações avançadas
- ▶ Combinação de DataFrames: `merge()`, `concat()`, `join()`
- ▶ Pivot tables e reshape
- ▶ Tratamento de duplicatas

Preparação:

- ▶ Revisar conceitos desta aula
- ▶ Praticar filtragem e seleção
- ▶ Experimentar com Titanic dataset
- ▶ Pensar em transformações que você gostaria de fazer

Dicas de Estudo

Para dominar Pandas:

1. Pratique com dados reais

- ▶ Kaggle datasets
- ▶ Dados públicos (IBGE, data.gov)
- ▶ Seus próprios dados

2. Documentação oficial

- ▶ <https://pandas.pydata.org/docs/>
- ▶ Guia do usuário muito completo
- ▶ Exemplos práticos

3. Cheat sheets

- ▶ DataCamp Pandas Cheat Sheet
- ▶ Pandas oficial cheat sheet

4. Explore gradualmente

- ▶ Comece simples (filtros básicos)
- ▶ Aumente complexidade aos poucos

Recursos Adicionais

Livros:

- ▶ "Python for Data Analysis- Wes McKinney (criador do Pandas)
- ▶ "Pandas Cookbook- Theodore Petrou

Tutoriais Online:

- ▶ Pandas oficial: 10 minutes to Pandas
- ▶ Real Python: Pandas tutorials
- ▶ Kaggle Learn: Pandas micro-course

Vídeos:

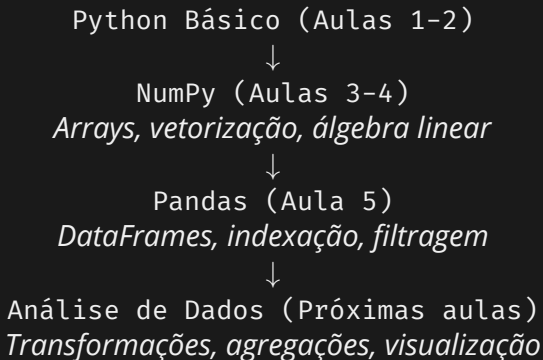
- ▶ Corey Schafer: Pandas Tutorial Series
- ▶ Keith Galli: Complete Python Pandas Data Science Tutorial

Prática:

- ▶ Kaggle: Titanic, House Prices, etc.
- ▶ DataCamp: Interactive exercises

Resumo: Do NumPy ao Pandas

Progressão do curso:



Exercício Prático

Tempo: 60 minutos

Entrega: via Moodle (notebook)

Tarefas:

1. (atualizado durante a aula)

Notebook: Disponível no Moodle

Obrigado!

Próxima aula: Manipulação de Dados com Pandas
Quinta-feira, 23/10

Dúvidas: via Moodle ou atendimento