

Programação para Ciência de Dados

Manipulação de Dados com Pandas

Arthur Casals

23 de Outubro de 2025

Agenda

- ▶ Transformações e Apply
- ▶ GroupBy e Agregações
- ▶ Combinação de DataFrames
- ▶ Reshape e Pivot

Recapitulação Aula 05

O que já sabemos:

- ▶ Series e DataFrames: estruturas fundamentais
- ▶ Criação: dicionários, listas, CSV
- ▶ Exploração: head(), info(), describe()
- ▶ Seleção: [], loc[], iloc[]
- ▶ Filtragem: boolean indexing, query()
- ▶ Ordenação: sort_values(), nlargest()

Hoje vamos aprofundar:

- ▶ Transformar dados de múltiplas formas
- ▶ Análise por grupos (split-apply-combine)
- ▶ Combinar múltiplos DataFrames
- ▶ Reestruturar dados (pivot, melt)

Bloco 1

Transformações e Apply

Transformações de Dados

Por que transformar dados?

- ▶ Criar novas features (feature engineering)
- ▶ Normalizar e padronizar valores
- ▶ Corrigir inconsistências
- ▶ Preparar dados para modelagem
- ▶ Extrair informações de colunas existentes

Métodos principais:

- ▶ **drop():** Remover linhas/colunas
- ▶ **Operações vetorizadas:** Mais rápidas e simples
- ▶ **apply():** Aplicar função a linhas/colunas
- ▶ **map():** Transformar Series elemento por elemento
- ▶ **applymap():** Aplicar função a cada célula (deprecated)
- ▶ **String methods:** Operações específicas para texto

DataFrame.drop()

Método para remover linhas ou colunas de um DataFrame

</> Python

```
1 DataFrame.drop(  
2     labels=None,  
3     axis=0,  
4     index=None,  
5     columns=None,  
6     level=None,  
7     inplace=False,  
8     errors='raise'  
9 )  
10 # Retorna: DataFrame ou None  
11 # Retorna um novo DataFrame sem os rótulos especificados (padrão)  
12 # Retorna None quando \texttt{inplace=True}
```

DataFrame.drop() - Parâmetros Básicos

| Parâmetro | Tipo/Valores | Descrição |
|----------------------|--------------------------|---|
| <code>labels</code> | único/lista | Rótulo(s) a remover (nomes de linhas ou colunas) |
| <code>axis</code> | 0/'index' ou 1/'columns' | Eixo para remover: 0=linhas, 1=colunas |
| <code>index</code> | único/lista | Alternativa: especificar rótulos de linhas diretamente |
| <code>columns</code> | único/lista | Alternativa: especificar rótulos de colunas diretamente |

💡 Nota Importante

Boa Prática: Use os parâmetros explícitos `index=` ou `columns=` para maior clareza ao invés de `labels` com `axis`.

DataFrame.drop() - Parâmetros Avançados

| Parâmetro | Tipo/Valores | Descrição |
|-----------|--------------------|---|
| level | int/nome | Para MultiIndex: especifica qual nível remover |
| inplace | bool | True: modifica o DataFrame original, retorna None False (padrão): retorna novo DataFrame |
| errors | 'raise' / 'ignore' | 'raise' (padrão): lança erro se rótulo não existe 'ignore': suprime erros |

💡 Nota Importante

Importante: Por padrão, `drop()` **não modifica** o DataFrame original. Sempre armazene o resultado ou use `inplace=True`.

Adicionar e Remover Colunas

Python

```
1 df = pd.read_csv('titanic.csv')
2 # Adicionar nova coluna (constante)
3 df['Status'] = 'Passenger'
4
5 # Adicionar baseada em outra coluna
6 df['AgeDouble'] = df['Age'] * 2
7
8 # Adicionar com condicao
9 df['IsChild'] = df['Age'] < 18
10
11 # Remover coluna
12 df = df.drop('AgeDouble', axis=1)
13
14 # Remover multiplas colunas
15 df = df.drop(['Status', 'IsChild'], axis=1)
```

Adicionar Colunas: Formas Alternativas

Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Usando assign (nao modifica original)
4 df_new = df.assign(
5     AgeGroup='Adult',
6     FareDouble=df['Fare'] * 2
7 )
8
9 # Multiplas colunas de uma vez
10 df[['Col1', 'Col2', 'Col3']] = [[1, 2, 3]] * len(df)
11
12 # Insert em posicao especifica
13 df.insert(2, 'NewCol', 0) # Posicao 2, nome, valor
14
```

Operações Vetorizadas: Preferência

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Operacoes aritmeticas (vetorizadas - RAPIDO)
4 df['FareEuros'] = df['Fare'] * 0.85
5 df['AgePlus10'] = df['Age'] + 10
6 df['FareSquared'] = df['Fare'] ** 2
7
8 # Operacoes entre colunas
9 df['FamilySize'] = df['SibSp'] + df['Parch'] + 1
10
11 # Operacoes condicionais vetorizadas
12 df['IsExpensive'] = df['Fare'] > 50
13
```

Operações Vetorizadas: Preferência

Nota Importante

Sempre prefira operações vetorizadas - são muito mais rápidas!

apply(), map(), applymap(): Diferenças

| Método | Aplica a | Retorna |
|------------|---|---------------------------|
| apply() | DataFrame/Series (linhas ou colunas) | Series/DataFrame |
| map() | Series (elemento por elemento) | Series |
| applymap() | DataFrame (cada célula) | DataFrame (deprecated) |

Quando usar:

- ▶ **apply():** Funções complexas em linhas/colunas
- ▶ **map():** Transformações elemento-a-elemento em Series
- ▶ **Vetorização:** Sempre que possível (mais rápido)
 - ▶ where(), select(), cut(), qcut(), acessores (.str, .dt), pd.Categorical...
 - ▶ Numba: decoradores @vectorize, @jit

apply(): Aplicar Função a Colunas

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Aplicar funcao a cada coluna (axis=0, padrao)
4 df_numeric = df[['Age', 'Fare']]
5 medias = df_numeric.apply(np.mean)
6
7 print(medias)
8 # Age      29.69
9 # Fare     32.20
```

apply(): Aplicar Função a Colunas (cont.)

</> Python

```
1 # Funcao customizada
2 def amplitude(col):
3     return col.max() - col.min()
4
5 amplitudes = df_numeric.apply(amplitude)
6
7 print(amplitudes)
8 # Age      79.5800
9 # Fare     512.3292
```

apply(): Aplicar Função a Linhas

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Aplicar funcao a cada linha (axis=1)
3 def calcular_familia(row):
4     return row['SibSp'] + row['Parch'] + 1
5
6 df['FamilySize'] = df.apply(calcular_familia, axis=1)
7
8 # Funcao inline (lambda)
9 df['FamilySize'] = df.apply(
10     lambda row: row['SibSp'] + row['Parch'] + 1,
11     axis=1
12 )
13
14 # MELHOR: vetorizado (evite apply se possivel)
15 df['FamilySize'] = df['SibSp'] + df['Parch'] + 1
```


apply(): Retornar Múltiplos Valores

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Funcao que retorna Series
3 def stats(col):
4     return pd.Series({
5         'min': col.min(),
6         'max': col.max(),
7         'mean': col.mean()
8     })
9
10 resultado = df[['Age', 'Fare']].apply(stats)
11 print(resultado)
```

| # | Age | Fare |
|--------|-------|--------|
| # min | 0.42 | 0.00 |
| # max | 80.00 | 512.33 |
| # mean | 29.69 | 32.20 |

apply(): Exemplo Complexo

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Extrair titulo do nome (Mr., Mrs., Miss., etc.)
3 def extrair_titulo(nome):
4     titulo = nome.split(',')[1].split('.')[0].strip()
5     return titulo
6
7 df['Titulo'] = df['Name'].apply(extrair_titulo)
8
9 print(df['Titulo'].value_counts().head())
10 # Mr          517
11 # Miss        182
12 # Mrs         125
13 # Master      40
14 # Dr           7
```

map(): Transformar Series

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Map com dicionario
4 mapa_classe = {1: 'First', 2: 'Second', 3: 'Third'}
5 df['ClasseName'] = df['Pclass'].map(mapa_classe)
6
7 # Map com funcao
8 df['SexLong'] = df['Sex'].map(
9     lambda x: 'Male' if x == 'male' else 'Female'
10 )
11
12 # Map com Series (como um dicionario)
13 medias_por_classe = df.groupby('Pclass')['Fare'].mean()
14 df['FareMedia'] = df['Pclass'].map(medias_por_classe)
```

map(): Transformar Series

</> Python

```
1 print(medias_por_classe)
2 # Pclass
3 # 1      84.154687
4 # 2      20.662183
5 # 3      13.675550
```

map() vs apply() em Series

Python

```
1 df = pd.read_csv('titanic.csv')
2 # map: apenas para Series
3 df['Sex'].map(lambda x: x.upper())
4 # apply: funciona para Series e DataFrame
5 df['Sex'].apply(lambda x: x.upper())
6
7 # Para Series simples, map e apply sao similares
8 # map é ligeiramente mais rapido para mapeamentos
9 # apply é mais flexivel (aceita funcoes complexas)
10 def processar(x):
11     # Logica complexa
12     return resultado
13
14 df['Sex'].apply(processar)
```

replace(): Substituição de Valores

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Substituir valores especificos
3 df['Embarked'] = df['Embarked'].replace('S', 'Southampton')
4
5 # Multiplas substituiçoes
6 df['Embarked'] = df['Embarked'].replace({
7     'S': 'Southampton',
8     'C': 'Cherbourg',
9     'Q': 'Queenstown'
10 })
11 # Substituir em todo o DataFrame
12 df = df.replace({'male': 'M', 'female': 'F'})
13
14 # Substituir valores faltantes
15 df = df.replace(np.nan, 0)
```

String Methods: str Accessor

Operações com strings no Pandas:

- ▶ Acessar via `.str`
- ▶ Métodos vetorizados (aplicados a toda Series)
- ▶ Similar a métodos de strings Python
- ▶ Lidam automaticamente com NaN

Principais categorias:

- ▶ **Case:** `lower()`, `upper()`, `capitalize()`
- ▶ **Limpeza:** `strip()`, `replace()`
- ▶ **Divisão:** `split()`, `extract()`
- ▶ **Busca:** `contains()`, `startswith()`, `endswith()`
- ▶ **Contagem:** `len()`, `count()`

String Methods: Case e Limpeza

Python

```
1 df = pd.read_csv('titanic.csv')
2 # Minusculas/Maiusculas
3 df['NameLower'] = df['Name'].str.lower()
4 df['NameUpper'] = df['Name'].str.upper()
5 df['NameTitle'] = df['Name'].str.title()
6
7 # Remover espaços
8 df['NameClean'] = df['Name'].str.strip()
9
10 # Substituir substring
11 df['NameMod'] = df['Name'].str.replace('Mr.', 'Mister')
12
13 # Remover caracteres especiais
14 df['NameAlpha'] = df['Name'].str.replace('[^a-zA-Z ]', '',
15                                           regex=True)
```


String Methods: Divisão e Extração

Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Dividir string
4 # Name: "Braund, Mr. Owen Harris"
5 sobrenome = df['Name'].str.split(',').str[0]
6 resto = df['Name'].str.split(',').str[1]
7
8 # Split com expand (retorna DataFrame)
9 nome_split = df['Name'].str.split(',', expand=True)
10 df['Sobrenome'] = nome_split[0]
11 df['Resto'] = nome_split[1]
12
13 # Extrair com regex
14 df['Titulo'] = df['Name'].str.extract(r'([^.]+\.)\.')
15
```

String Methods: Busca e Verificação

Python

```
1 df = pd.read_csv('titanic.csv')
2 # Verificar se contem substring
3 tem_mr = df['Name'].str.contains('Mr.')
4 print(tem_mr.sum()) # Quantos "Mr."
5 # Comeca com
6 comeca_b = df['Name'].str.startswith('B')
7 # Termina com
8 termina_s = df['Name'].str.endswith('s')
9 # Case insensitive
10 tem_john = df['Name'].str.contains('john', case=False)
11 # Contar ocorrencias
12 n_espacos = df['Name'].str.count(' ')
```

String Methods: Tamanho e Padding

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Tamanho da string
3 df['NameLength'] = df['Name'].str.len()
4
5 # Slice (fatiar)
6 primeiras_5 = df['Name'].str[:5]
7 ultimas_3 = df['Name'].str[-3:]
8
9 # Padding (preencher com caracteres)
10 df['IdPadded'] = df['PassengerId'].astype(str).str.zfill(5)
11 # 1 -> 00001, 42 -> 00042
12
13 # Concatenar strings
14 df['Info'] = df['Name'] + ' (' + df['Sex'] + ')'
```

Exemplo Prático: Extrair Título

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Extrair título do nome
3 # "Braund, Mr. Owen Harris" -> "Mr"
4 df['Titulo'] = df['Name'].str.extract(r', ([^.]+)\.')
5 # Limpar espaços
6 df['Titulo'] = df['Titulo'].str.strip()
7 # Ver distribuicao
8 print(df['Titulo'].value_counts())
9 # Agrupar titulos raros
10 titulos_raros = ['Don', 'Rev', 'Dr', 'Mme', 'Ms',
11                  'Major', 'Lady', 'Sir', 'Mlle',
12                  'Col', 'Capt', 'Countess', 'Jonkheer']
13 df['TituloGroup'] = df['Titulo'].replace(titulos_raros, 'Rare')
```

Conversão de Tipos

Por que converter tipos?

- ▶ Otimizar memória (int64 → int32)
- ▶ Corrigir leitura incorreta de CSV
- ▶ Preparar dados para operações específicas
- ▶ Converter categóricas para economia de memória

Métodos principais:

- ▶ **astype():** Conversão forçada
- ▶ **to_numeric():** Strings → números (com error handling)
- ▶ **to_datetime():** Strings → datas
- ▶ **to_timedelta():** Strings → intervalos de tempo
- ▶ **astype('category'):** Otimizar categóricas

astype(): Conversão Básica

Python

```
1 df = pd.read_csv('titanic.csv')
2 # Converter para diferente tipo numerico
3 df['Pclass'] = df['Pclass'].astype('int32')
4 df['Age'] = df['Age'].astype('float32')
5 # Converter para string
6 df['PassengerId'] = df['PassengerId'].astype(str)
7 # Converter booleano
8 df['Survived'] = df['Survived'].astype(bool)
```

astype(): Conversão Básica (cont.)

</> Python

```
1 # Verificar tipos
2 print(df.dtypes)
3 # PassengerId      object
4 # Survived         bool
5 # Pclass           int32
6 # ...
7
8 # Memoria usada
9 print(df.memory_usage(deep=True))
10 # Index            132
11 # PassengerId      46224
12 # Survived         891
13 # ...
```

to_numeric(): Conversão Segura

</> Python

```
1 # Dados com valores invalidos
2 s = pd.Series(['1', '2', '3.5', 'four', '5'])
3 # astype vai dar erro: s.astype(float) => ValueError!
4 # to_numeric com error handling
5 s_numeric = pd.to_numeric(s, errors='coerce')
6 print(s_numeric)
7 # 0      1.0
8 # 1      2.0
9 # 2      3.5
10 # 3      NaN    <- invalido virou NaN
11 # 4      5.0
12 # dtype: float64
```


to_numeric(): Conversão Segura

</> Python

```
1 # Ignorar erros (manter original)
2 s_numeric = pd.to_numeric(s, errors='ignore')
3 print(s_numeric)
4 # 0      1.0
5 # 1      2.0
6 # 2      3.5
7 # 3     four
8 # 4      5.0
9 # dtype: object  <- tipo agora é objeto!
```

Categorical: Otimização de Memória

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Memória antes
3 print(f"Memória Sex: {df['Sex'].memory_usage(deep=True)} bytes")
4 # Memória Sex: 47983 bytes
5 # Converter para categorical
6 df['Sex'] = df['Sex'].astype('category')
7 df['Embarked'] = df['Embarked'].astype('category')
8 df['Pclass'] = df['Pclass'].astype('category')
9 # Memória depois
10 print(f"Memória Sex: {df['Sex'].memory_usage(deep=True)} bytes")
11 # Memória Sex: 1239 bytes
12 print(df['Sex'].cat.categories)
13 # Index(['female', 'male'], dtype='object')
```

Binning: Discretização de Variáveis

O que é binning?

- ▶ Converter variável contínua em categórica
- ▶ Agrupar valores em "bins" (faixas)
- ▶ Facilita análise e visualização
- ▶ Reduz ruído e outliers

Dois métodos principais:

- ▶ **pd.cut()**: Bins de largura igual
- ▶ **pd.qcut()**: Bins com frequência igual (quantis)

💡 Nota Importante

Escolha baseada no objetivo: distribuição uniforme (qcut) ou intervalos fixos (cut)

pd.cut(): Bins de Largura Igual

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Criar faixas de idade com largura igual
3 df['AgeGroup'] = pd.cut(df['Age'], bins=4) # 4 bins automáticos
4 # Criar faixas de idade pré-definidas
5 df['AgeGroup'] = pd.cut(
6     df['Age'],
7     bins=[0, 18, 40, 60, 100], # (0, 18], (18, 40], (40, 60], (60, 100]
8     labels=['Crianca', 'Jovem', 'Adulto', 'Idoso']
9 )
10 print(df['AgeGroup'].value_counts())
11 # Jovem      357
12 # Adulto     241
13 # Crianca    83
14 # Idoso      33
15
```

pd.qcut(): Bins por Quantis

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Criar faixas com frequência similar (quartis)
3 df['FareQuartile'] = pd.qcut(
4     df['Fare'],
5     q=4,
6     labels=['Q1-Baixo', 'Q2-Medio', 'Q3-Alto', 'Q4-Premium']
7 )
8 print(df['FareQuartile'].value_counts())
9 # Q1-Baixo      223
10 # Q2-Medio      223
11 # Q3-Alto       223
12 # Q4-Premium    222
13 # Ver os limites criados
14 print(df['FareQuartile'].cat.categories)
15 # Index(['Q1-Baixo', 'Q2-Medio', 'Q3-Alto', 'Q4-Premium'], dtype='object')
```

Binning: Comparação cut vs qcut

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # cut: intervalos fixos
3 df['Fare_Cut'] = pd.cut(df['Fare'], bins=4)
4 print("Distribuição com cut():")
5 print(df['Fare_Cut'].value_counts().sort_index())
6 # (-0.513, 128.08]      871  <- DESBALANCEADO!
7 # (128.08, 256.16]      12
8 # (256.16, 384.24]       6
9 # (384.24, 512.33]       2
10 # qcut: frequências iguais
11 df['Fare_Qcut'] = pd.qcut(df['Fare'], q=4)
12 print("\nDistribuição com qcut():")
13 print(df['Fare_Qcut'].value_counts().sort_index())
14 # Todos com ~223 elementos (BALANCEADO!)
```

Binning: Opções Avançadas

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Incluir bins abertos à direita ou esquerda
3 df['Age_Right'] = pd.cut(df['Age'], bins=[0, 18, 60, 100],
4                           right=True)    # (0, 18]
5 df['Age_Left'] = pd.cut(df['Age'], bins=[0, 18, 60, 100],
6                           right=False)   # [0, 18)
7 # Retornar os bins como inteiros
8 df['Age_Bin'] = pd.cut(df['Age'], bins=5, labels=False)
9 # 0, 1, 2, 3, 4
10 # Incluir valores fora dos limites
11 df['Age_Extended'] = pd.cut(df['Age'], bins=[0, 18, 60],
12                             labels=['Jovem', 'Adulto'])
13 # NaN para idade > 60
```

Quando Usar cut vs qcut

Use `pd.cut()` quando:

- ▶ Limites têm significado real
- ▶ Faixas fixas são importantes
- ▶ Exemplo: Idade (0-18, 18-65, 65+)
- ▶ Exemplo: Notas (0-49, 50-69, 70-100)
- ▶ Exemplo: IMC (abaixo, normal, sobrepeso)

Use `pd.qcut()` quando:

- ▶ Quer distribuição balanceada
- ▶ Valores extremos (outliers)
- ▶ Machine learning (features balanceadas)
- ▶ Exemplo: Segmentação de clientes
- ▶ Exemplo: Ranking percentual

💡 Nota Importante

`qcut` é útil para evitar bins vazios com dados assimétricos

Exemplo Prático: Pipeline de Transformação

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 print(f"Features: {df.shape[1]} colunas")
3 # Features: 12 colunas
4 # 1. Criar FamilySize
5 df['FamilySize'] = df['SibSp'] + df['Parch'] + 1
6 # 2. Extrair título
7 df['Titulo'] = df['Name'].str.extract(r', ([^.]+)\.')
8 # 3. Criar IsAlone
9 df['IsAlone'] = (df['FamilySize'] == 1).astype(int)
10 # 4. Categorizar idade
11 df['AgeGroup'] = pd.cut(df['Age'],
12                          bins=[0, 12, 18, 60, 100],
13                          labels=['Child', 'Teen', 'Adult', 'Senior'])
```

Exemplo Prático: Pipeline (cont.)

</> Python

```
1 # 5. Converter para categorical
2 for col in ['Sex', 'Embarked', 'Pclass', 'Titulo']:
3     df[col] = df[col].astype('category')
4
5 # 6. Normalizar tarifa
6 df['FareNorm'] = (df['Fare'] - df['Fare'].min()) / \
7                 (df['Fare'].max() - df['Fare'].min())
8
9 # 7. Preencher idade faltante com mediana
10 df['Age'].fillna(df['Age'].median(), inplace=True)
11
12 print(f"Features: {df.shape[1]} colunas")
13 # Features: 17 colunas
14 print(df.head())
15 # Check new columns
```

Window Functions: Operações em Janelas

O que são Window Functions?

- ▶ Operações que consideram valores "vizinhos"
- ▶ Não agregam (mantêm número de linhas)
- ▶ Úteis para comparações temporais/ordenadas
- ▶ Similar ao SQL window functions

Principais funções:

- ▶ **shift()**: Valores anteriores/posteriores
- ▶ **diff()**: Diferença entre valores consecutivos
- ▶ **pct_change()**: Variação percentual
- ▶ **rolling()**: Médias móveis
- ▶ **rank()**: Ranking dentro de grupos
- ▶ **cumsum()/cumcount()**: Acumulados

shift(): Acessar Valores Anteriores

</> Python

```
1 vendas = pd.DataFrame({ # Dados de vendas mensais
2     'mes': ['Jan', 'Fev', 'Mar', 'Abr', 'Mai'],
3     'valor': [100, 120, 115, 140, 150]
4 })
5 vendas['mes_anterior'] = vendas['valor'].shift(1)
6 vendas['mes_seguinte'] = vendas['valor'].shift(-1)
7 print(vendas)
8 #   mes  valor  mes_anterior  mes_seguinte
9 #  0   Jan    100          NaN         120.0
10 #  1   Fev    120         100.0         115.0
11 #  2   Mar    115         120.0         140.0
12 #  3   Abr    140         115.0         150.0
13 #  4   Mai    150         140.0          NaN
```

diff() e pct_change(): Variações

Python

```
1 vendas = pd.DataFrame({
2     'mes': ['Jan', 'Fev', 'Mar', 'Abr'],
3     'valor': [100, 120, 115, 140]
4 })
5 # Diferença absoluta
6 vendas['variacao'] = vendas['valor'].diff()
7 # Variação percentual
8 vendas['pct_variacao'] = vendas['valor'].pct_change()
9 print(vendas)
10 #   mes  valor  variacao  pct_variacao
11 # 0  Jan    100      NaN             NaN
12 # 1  Fev    120    20.0      0.200000    # +20%
13 # 2  Mar    115   -5.0     -0.041667    # -4.17%
14 # 3  Abr    140   25.0      0.217391    # +21.74%
```

rank(): Ranking de Valores

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Ranking geral de tarifa (maior = 1)
3 df['FareRank'] = df['Fare'].rank(ascending=False, method='dense')
4 # Ranking dentro de cada classe
5 df['FareRankClass'] = df.groupby('Pclass')['Fare'].rank(
6     ascending=False,
7     method='dense'
8 )
9 # Métodos de ranking
10 # 'average': média das posições (padrão)
11 # 'min': menor posição
12 # 'max': maior posição
13 # 'first': ordem de aparição
14 # 'dense': sem gaps (1, 2, 3, ...)
```

rank(): Ranking de Valores (cont.)

</> Python

```
1 print(df[['Name', 'Pclass', 'Fare', 'FareRankClass']].head(10))
2 #      [Name...] Pclass      Fare  FareRankClass
3 # 0              3    7.2500          102.0
4 # 1              1   71.2833           40.0
5 # 2              3    7.9250           79.0
6 # 3              1   53.1000           55.0
7 # ...
```

cumsum() e cumcount(): Acumulados

</> Python

```
1 vendas = pd.DataFrame({
2     'dia': range(1, 6),
3     'valor': [100, 120, 115, 140, 150]
4 })
5 # Soma acumulada
6 vendas['acumulado'] = vendas['valor'].cumsum()
7 # Contagem acumulada (útil com groupby)
8 vendas['count'] = vendas['valor'].cumcount()
9 print(vendas)
10 #    dia  valor  acumulado  count
11 # 0     1    100         100     0
12 # 1     2    120         220     1
13 # 2     3    115         335     2
14 # 3     4    140         475     3
15 # 4     5    150         625     4
```


Window Functions com GroupBy

Python

```
1 df = pd.read_csv('titanic.csv')
2 # Ranking de idade dentro de cada classe
3 df['AgeRank'] = df.groupby('Pclass')['Age'].rank(
4     ascending=False
5 )
6 # Diferença da idade média da classe
7 df['AgeDiff'] = df.groupby('Pclass')['Age'].transform(
8     lambda x: x - x.mean()
9 )
10 # Acumulado de tarifa por classe (ordem alfabética de nome)
11 df = df.sort_values(['Pclass', 'Name'])
12 df['FareCumsum'] = df.groupby('Pclass')['Fare'].cumsum()
```

Window Functions com GroupBy (cont.)

</> Python

```
1 print(df[['Pclass', 'Age', 'AgeRank', 'AgeDiff']].head())
2 #      Pclass    Age  AgeRank  AgeDiff
3 # 730         1  29.00    133.0   -9.233441
4 # 305         1   0.92    186.0  -37.313441
5 # 297         1   2.00    185.0 -36.233441
6 # 498         1  25.00    146.0 -13.233441
7 # 460         1  48.00     52.0   9.766559
```

Exemplo Prático: Análise de Crescimento

</> Python

```
1 # Vendas por mês
2 vendas_mensais = pd.DataFrame({
3     'mes': pd.date_range('2024-01', periods=6, freq='MS'),
4     'vendas': [1000, 1200, 1150, 1400, 1500, 1350]
5 })
6 # Pipeline completo de análise temporal
7 vendas_mensais['mes_anterior'] = vendas_mensais['vendas'].shift(1)
8 vendas_mensais['crescimento'] = vendas_mensais['vendas'].diff()
9 vendas_mensais['crescimento_pct'] = vendas_mensais['vendas'].pct_change()
10 vendas_mensais['acumulado'] = vendas_mensais['vendas'].cumsum()
11 vendas_mensais['media_movel_3'] = vendas_mensais['vendas'].rolling(
12     window=3
13 ).mean()
```

Exemplo Prático: Análise de Crescimento (cont.)

</> Python

```
1 print(vendas_mensais)
2 #           mes  vendas  mes_anterior  crescimento  crescimento_pct  \
3 # 0 2024-01-01    1000             NaN             NaN             NaN
4 # 1 2024-02-01    1200          1000.0           200.0           0.200000
5 # 2 2024-03-01    1150          1200.0           -50.0          -0.041667
6 # ...
7 #   acumulado  media_movel_3
8 # 0         1000             NaN
9 # 1         2200             NaN
10 # 2         3350          1116.666667
11 # ...
```

Window Functions: Casos de Uso

Cenários comuns:

1. Comparações temporais

- ▶ Vendas mês-a-mês
- ▶ Crescimento ano-a-ano (YoY)

2. Rankings

- ▶ Top N produtos por categoria
- ▶ Classificação de clientes

3. Métricas acumuladas

- ▶ Revenue acumulado
- ▶ Contagem progressiva de eventos

4. Médias móveis

- ▶ Suavização de tendências
- ▶ Detecção de anomalias

Bloco 2

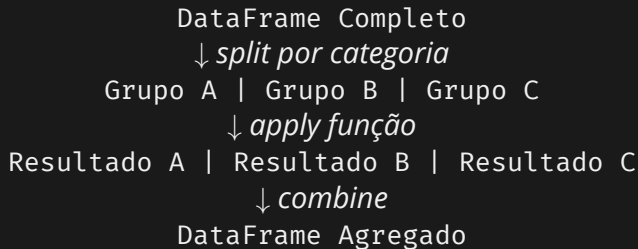
GroupBy e Agregações

GroupBy: Split-Apply-Combine

Conceito:

1. **Split:** Dividir dados em grupos
2. **Apply:** Aplicar função a cada grupo
3. **Combine:** Combinar resultados

Visual:



GroupBy: Quando Usar

Perguntas que GroupBy responde:

- ▶ Qual a **média** por categoria?
- ▶ Quantos elementos em cada grupo?
- ▶ Qual o **máximo/mínimo** por grupo?
- ▶ Como os grupos se **comparam**?

Exemplos práticos:

- ▶ Taxa de sobrevivência por classe (Titanic)
- ▶ Vendas totais por região
- ▶ Salário médio por departamento
- ▶ Nota média por turma
- ▶ Métricas por segmento de clientes

GroupBy: Sintaxe Básica

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Agrupar por uma coluna
3 grouped = df.groupby('Pclass')
4 print(type(grouped)) # DataFrameGroupBy object
5 # Ver grupos
6 print(grouped.groups.keys())
7 # dict_keys([1, 2, 3])
8 # Numero de elementos em cada grupo
9 print(grouped.size())
10 # Pclass
11 # 1      216
12 # 2      184
13 # 3      491
```

GroupBy: Agregações Simples

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Agrupar e calcular media
4 media_por_classe = df.groupby('Pclass')['Age'].mean()
5 print(media_por_classe)
6 # Pclass
7 # 1      38.23
8 # 2      29.88
9 # 3      25.14
10
11 # Múltiplas estatísticas
12 stats = df.groupby('Pclass')['Age'].agg(['mean', 'std',
13                                           'min', 'max'])
14 print(stats)
15 # Colunas: mean, std, min, max
```

GroupBy: Múltiplas Colunas

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Agrupar por multiplas colunas
3 grouped = df.groupby(['Pclass', 'Sex'])
4
5 # Taxa de sobrevivencia por classe e sexo
6 sobrev = grouped['Survived'].mean()
7 print(sobrev)
8 # Pclass  Sex
9 # 1      female    0.968
10 #      male      0.369
11 # 2      female    0.921
12 #      male      0.157
13 # 3      female    0.500
14 #      male      0.135
```

GroupBy: Agregações em Múltiplas Colunas

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Agregar multiplas colunas
3 resultado = df.groupby('Pclass')[['Age', 'Fare']].mean()
4 print(resultado)
5 #           Age           Fare
6 # Pclass
7 # 1      38.23      84.15
8 # 2      29.88      20.66
9 # 3      25.14      13.68
10 # Diferentes funcoes para diferentes colunas
11 resultado = df.groupby('Pclass').agg({
12     'Age': 'mean',
13     'Fare': 'sum',
14     'Survived': 'count'
15 })
```

agg(): Agregações Personalizadas

Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Multiplas funcoes para uma coluna
4 resultado = df.groupby('Pclass')['Age'].agg([
5     'count',
6     'mean',
7     'std',
8     'min',
9     ('q25', lambda x: x.quantile(0.25)),
10    ('q75', lambda x: x.quantile(0.75)),
11    'max'
12 ])
```

agg(): Agregações Personalizadas (cont.)

</> Python

```
1 print(resultado)
2 #              count          mean          std          min          q25          q75          max
3 # Pclass
4 # 1             186   38.233441   14.802856   0.92    27.0    49.0    80.0
5 # 2             173   29.877630   14.001077   0.67    23.0    36.0    70.0
6 # 3             355   25.140620   12.495398   0.42    18.0    32.0    74.0
7
8
```

agg(): Funções Customizadas

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Funcao customizada
3 def amplitude(x):
4     return x.max() - x.min()
5 # Usar funcao customizada
6 resultado = df.groupby('Pclass')['Age'].agg([
7     'mean',
8     amplitude,
9     ('range', amplitude) # Com nome customizado
10 ])
11 # Funcoes lambda
12 resultado = df.groupby('Pclass')['Fare'].agg([
13     ('soma', 'sum'),
14     ('dobro_media', lambda x: x.mean() * 2)
15 ])
```

agg(): Diferentes Funções por Coluna

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Dicionario: coluna -> funcoes
3 resultado = df.groupby('Pclass').agg({
4     'Age': ['mean', 'std'],
5     'Fare': ['sum', 'mean', 'max'],
6     'Survived': ['sum', 'mean'],
7     'PassengerId': 'count'
8 })
```


agg(): Diferentes Funções por Coluna (cont.)

</> Python

```
1 print(resultado)
2 # MultiIndex em colunas - rótulos são tuplas: ('Age', 'mean')
3 #           Age           Fare           Survived
4 #           mean      std      sum      mean      max      sum
5 # Pclass
6 # 1      38.233441  14.802856  18177.4125  84.154687  512.3292  136
7 # 2      29.877630  14.001077   3801.8417  20.662183   73.5000   87
8 # 3      25.140620  12.495398   6714.6951  13.675550   69.5500  119
9 # \
10 #           PassengerId
11 #           mean      count
12 # Pclass
13 # 1      0.629630      216
14 # 2      0.472826      184
15 # 3      0.242363      491
```

agg(): Diferentes Funções por Coluna (cont.)

</> Python

```
1 # Achatar colunas: ('Age', 'mean') => Age_mean
2 resultado.columns = ['_'.join(col) for col in resultado.columns]
3 print(resultado)
4 #           Age_mean      Age_std      Fare_sum      Fare_mean      Fare_max [...]
5 # Pclass
6 # 1           38.233441    14.802856    18177.4125     84.154687     512.3292
7 # 2           29.877630    14.001077     3801.8417     20.662183      73.5000
8 # 3           25.140620    12.495398     6714.6951     13.675550     69.5500
```

Named Aggregations: Sintaxe Moderna

Problema com agg() tradicional:

- ▶ Colunas com MultiIndex difíceis de usar
- ▶ Precisa renomear manualmente depois
- ▶ Código menos legível

Solução: Named Aggregations

- ▶ Nomear agregações diretamente
- ▶ Sintaxe clara: nome=(coluna, função)
- ▶ Sem MultiIndex nos resultados
- ▶ Mais fácil de documentar

💡 Nota Importante

Recomendado para código de produção - muito mais legível!

Named Aggregations: Sintaxe

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # ANTES: agg tradicional (difícil de ler)
4 resultado_antigo = df.groupby('Pclass').agg({
5     'Age': ['mean', 'std'],
6     'Fare': ['sum', 'mean'],
7     'Survived': 'sum'
8 })
9 # Colunas: ('Age', 'mean'), ('Age', 'std'), ...
10 resultado_antigo.columns = ['_'.join(col) for col in resultado_antigo.
    columns]
```

Named Aggregations: Sintaxe (cont.)

</> Python

```
1 # AGORA: Named aggregations (limpo!)
2 resultado_novo = df.groupby('Pclass').agg(
3     idade_media=('Age', 'mean'),
4     idade_desvio=('Age', 'std'),
5     receita_total=('Fare', 'sum'),
6     ticket_medio=('Fare', 'mean'),
7     sobreviventes=('Survived', 'sum')
8 )
```

Named Aggregations: Sintaxe (cont.)

</> Python

```
1 print(resultado_novo)
2 #          idade_media  idade_desvio  receita_total  ticket_medio
   sobreviventes
3 # Pclass
4 # 1          38.233441      14.802856      18177.4125      84.154687
   136
5 # 2          29.877630      14.001077       3801.8417      20.662183
   87
6 # 3          25.140620      12.495398       6714.6951      13.675550
   119
```

Named Aggregations: Múltiplos Grupos

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Agregações nomeadas com múltiplos grupos
4 analyse = df.groupby(['Pclass', 'Sex']).agg(
5     total_passageiros=('PassengerId', 'count'),
6     total_sobreviventes=('Survived', 'sum'),
7     taxa_sobrevivencia=('Survived', 'mean'),
8     idade_media=('Age', 'mean'),
9     tarifa_mediana=('Fare', 'median'),
10    tarifa_maxima=('Fare', 'max')
11 ).round(2)
```

Named Aggregations: Múltiplos Grupos (cont.)

</> Python

```
1 print(analise) # Colunas com nomes claros, sem MultiIndex!
2 #              total_passageiros  total_sobreviventes  [...]
3 # Pclass Sex
4 # 1      female                94                91
5 #      male                  122                45
6 # 2      female                76                70
7 #      male                  108                17
8 # 3      female                144                72
9 #      male                  347                47
```


Named Aggregations: Funções Customizadas

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Funções customizadas com named agg
3 def amplitude(x):
4     return x.max() - x.min()
5 def percentil_90(x):
6     return x.quantile(0.9)
7 resultado = df.groupby('Pclass').agg(
8     idade_min=('Age', 'min'),
9     idade_max=('Age', 'max'),
10    idade_amplitude=('Age', amplitude),
11    tarifa_p90=('Fare', percentil_90),
12    tarifa_total=('Fare', 'sum')
13 )
```

Named Aggregations: Com Lambda

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Lambda functions com named aggregations
3 resultado = df.groupby('Pclass').agg(
4     # Funções built-in
5     total=('PassengerId', 'count'),
6     idade_media=('Age', 'mean'),
7     # Lambda para cálculos customizados
8     idade_cv=('Age', lambda x: x.std() / x.mean()), # Coef. variação
9     tarifa_iqr=('Fare', lambda x: x.quantile(0.75) - x.quantile(0.25)),
10    sobrev_pct=('Survived', lambda x: x.mean() * 100)
11 ).round(2)
```

Named Aggregations: Vantagens

Por que usar Named Aggregations:

1. Legibilidade

- ▶ Código auto-documentado
- ▶ Nomes significativos para cada métrica

2. Sem pós-processamento

- ▶ Não precisa renomear colunas depois
- ▶ Sem MultiIndex para "achatar"

3. Fácil manutenção

- ▶ Adicionar/remover métricas é simples
- ▶ Fica claro o que cada agregação faz

4. Compatível com pipelines

- ▶ Integra bem com method chaining
- ▶ Pronto para usar em visualizações

Filtrar Grupos

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Criar grupos por idade
3 df['AgeGroup'] = pd.cut(df['Age'], bins=[0, 18, 60, 100],
4                           labels=['Child', 'Adult', 'Senior'])
5 # Filtrar grupos com mais de 100 pessoas
6 grupos_grandes = df.groupby('AgeGroup').filter(
7     lambda x: len(x) > 100
8 )
9 print(f"Original: {len(df)}")
10 print(f"Filtrado: {len(grupos_grandes)}")
11 # Apenas grupos onde sobrevivencia > 50%
12 sobrev_alta = df.groupby('Pclass').filter(
13     lambda x: x['Survived'].mean() > 0.5
14 )
```

transform() vs aggregate()

aggregate():

- ▶ Reduz grupos a valores únicos
- ▶ Retorna resultado menor
- ▶ Uma linha por grupo
- ▶ Ex: média, soma, count

Resultado:

- ▶ Shape menor que original
- ▶ Índice = grupos

transform():

- ▶ Mantém shape original
- ▶ Retorna valor para cada linha
- ▶ Broadcasting de agregação
- ▶ Ex: normalização por grupo

Resultado:

- ▶ Shape igual ao original
- ▶ Pode adicionar como coluna

transform(): Exemplo

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Calcular media de idade por classe
3 media_classe = df.groupby('Pclass')['Age'].transform('mean')
4
5 # Agora media_classe tem o mesmo tamanho que df!
6 print(len(media_classe)) # 891
7 print(len(df))          # 891
8
9 # Adicionar como nova coluna
10 df['AgeMeanClass'] = media_classe
11
12 # Normalizar idade dentro de cada classe
13 df['AgeNormClass'] = df.groupby('Pclass')['Age'].transform(
14     lambda x: (x - x.mean()) / x.std()
15 )
```

transform(): Casos de Uso

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # 1. Preencher NaN com media do grupo
3 df['Age'] = df.groupby('Pclass')['Age'].transform(
4     lambda x: x.fillna(x.mean())
5 )
6
7 # 2. Ranking dentro do grupo
8 df['FareRank'] = df.groupby('Pclass')['Fare'].transform(
9     lambda x: x.rank(ascending=False)
10 )
11
12 # 3. Diferença da media do grupo
13 df['FareDiff'] = df.groupby('Pclass')['Fare'].transform(
14     lambda x: x - x.mean()
15 )
```

apply(): Mais Flexível que agg/transform

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # apply pode retornar Series, DataFrame ou escalar
3 def top_fares(group):
4     return group.nlargest(3, 'Fare')[['Name', 'Fare']]
5 # Top 3 tarifas por classe
6 resultado = df.groupby('Pclass').apply(top_fares)
7 print(resultado)
8 # Funcao customizada complexa
9 def stats_customizado(group):
10     return pd.Series({
11         'count': len(group),
12         'survived_pct': group['Survived'].mean(),
13         'avg_age': group['Age'].mean()
14     })
15 resultado = df.groupby('Pclass').apply(stats_customizado)
```


Iteração sobre Grupos

Python

```
1 df = pd.read_csv('titanic.csv')
2 # Iterar sobre grupos
3 for nome_grupo, dados_grupo in df.groupby('Pclass'):
4     print(f"\n=== Classe {nome_grupo} ===")
5     print(f"Passageiros: {len(dados_grupo)}")
6     print(f"Sobreviventes: {dados_grupo['Survived'].sum()}")
7     print(f"Taxa: {dados_grupo['Survived'].mean():.2%}")
8 # === Classe 1 ===
9 # Passageiros: 216
10 # Sobreviventes: 136
11 # Taxa: 62.96%
12 grouped = df.groupby('Pclass')
13 primeira_classe = grouped.get_group(1)
14 print(primeira_classe.head())
15 # Visualiza só primeira classe
```

Exemplo Prático: Análise por Grupo

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Análise completa por classe e sexo
3 analise = df.groupby(['Pclass', 'Sex']).agg({
4     'PassengerId': 'count',
5     'Survived': ['sum', 'mean'],
6     'Age': ['mean', 'median'],
7     'Fare': ['mean', 'max']
8 })
9 # Renomear colunas
10 analise.columns = ['_'.join(col) for col in analise.columns]
11 analise = analise.rename(columns={
12     'PassengerId_count': 'Total',
13     'Survived_sum': 'Sobreviventes',
14     'Survived_mean': 'Taxa_Sobrev'
15 })
```

Exemplo Prático: Ranking por Grupo

Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Ranking de tarifa dentro de cada classe
4 df['FareRank'] = df.groupby('Pclass')['Fare'].rank(
5     ascending=False,
6     method='dense'
7 )
8
9 # Top 5 tarifas de cada classe
10 top5_por_classe = df[df['FareRank'] <= 5].sort_values(
11     ['Pclass', 'FareRank']
12 )
```

Bloco 3

Combinação de DataFrames

Por que Combinar DataFrames?

Cenários comuns:

- ▶ Dados em múltiplos arquivos
- ▶ Informações complementares em tabelas separadas
- ▶ Junção de datasets de diferentes fontes
- ▶ Adicionar features de lookup tables
- ▶ Consolidar dados temporais

Métodos principais:

- ▶ **concat():** Empilhar DataFrames (vertical/horizontal)
- ▶ **merge():** JOIN estilo SQL (por chaves)
- ▶ **join():** Merge pelo índice
- ▶ **append():** Adicionar linhas (deprecated, use concat)

concat() vs merge()

concat():

- ▶ **Empilhamento** de dados
- ▶ Não requer chave comum
- ▶ Baseado em posição/índice
- ▶ Útil para dados similares

Casos de uso:

- ▶ Múltiplos CSVs com mesma estrutura
- ▶ Dados de períodos diferentes
- ▶ Adicionar colunas

merge():

- ▶ **JOIN** estilo SQL
- ▶ Requer chave(s) comum(ns)
- ▶ Baseado em valores
- ▶ Útil para relacionamentos

Casos de uso:

- ▶ Juntar tabelas relacionadas
- ▶ Adicionar info de lookup
- ▶ Relacionamentos N:M

concat(): Empilhamento Vertical

</> Python

```
1 # Dois DataFrames com mesma estrutura
2 df1 = pd.DataFrame({
3     'A': [1, 2, 3],
4     'B': [4, 5, 6]
5 })
6
7 df2 = pd.DataFrame({
8     'A': [7, 8, 9],
9     'B': [10, 11, 12]
10 })
11
```

concat(): Empilhamento Vertical (cont.)

</> Python

```
1 # Concatenar verticalmente (empilhar linhas)
2 resultado = pd.concat([df1, df2])
3 print(resultado)
4 #      A      B
5 # 0    1      4
6 # 1    2      5
7 # 2    3      6
8 # 0    7     10  <- indice repetido!
9 # 1    8     11
10 # 2    9     12
11
```


concat(): Resetar Índice

</> Python

```
1 df1 = pd.DataFrame({'A': [1, 2], 'B': [3, 4]})
2 df2 = pd.DataFrame({'A': [5, 6], 'B': [7, 8]})
3
4 # Sem resetar indice
5 resultado = pd.concat([df1, df2])
6 print(resultado.index) # [0, 1, 0, 1]
7
8 # Ignorar indice original
9 resultado = pd.concat([df1, df2], ignore_index=True)
10 print(resultado.index) # [0, 1, 2, 3]
11
12 # Alternativa: resetar depois
13 resultado = pd.concat([df1, df2]).reset_index(drop=True)
14
```

concat(): Empilhamento Horizontal

</> Python

```
1 df1 = pd.DataFrame({
2     'A': [1, 2, 3],
3     'B': [4, 5, 6]
4 })
5 df2 = pd.DataFrame({
6     'C': [7, 8, 9],
7     'D': [10, 11, 12]
8 })
9 # Concatenar horizontalmente (adicionar colunas)
10 resultado = pd.concat([df1, df2], axis=1)
11 print(resultado)
12 #      A  B  C  D
13 # 0    1  4  7 10
14 # 1    2  5  8 11
15 # 2    3  6  9 12
```

concat(): Lidando com Colunas Diferentes

</> Python

```
1 df1 = pd.DataFrame({'A': [1, 2], 'B': [3, 4]})
2 df2 = pd.DataFrame({'B': [5, 6], 'C': [7, 8]})
3 # Por padrao: outer join (todas as colunas)
4 resultado = pd.concat([df1, df2])
5 print(resultado)
6 #      A      B      C
7 # 0  1.0  3.0  NaN
8 # 1  2.0  4.0  NaN
9 # 0  NaN  5.0  7.0
10 # 1  NaN  6.0  8.0
11
12 # Inner join (apenas colunas comuns)
13 resultado = pd.concat([df1, df2], join='inner')
14 print(resultado) # Apenas coluna B
```

concat(): Múltiplos DataFrames

</> Python

```
1 # Lista de DataFrames
2 dfs = []
3 for i in range(5): # Cria 5 DFs
4     df = pd.DataFrame({
5         'id': range(i*10, (i+1)*10),
6         'value': np.random.randn(10)
7     })
8     dfs.append(df)
9 # Concatenar todos de uma vez
10 resultado = pd.concat(dfs, ignore_index=True) # Cria índice novo
11 print(f"Total de linhas: {len(resultado)}") # 50
12 # Com identificador de origem
13 resultado = pd.concat(dfs, keys=['df1', 'df2', 'df3', 'df4', 'df5'])
14 print(resultado.index) # MultiIndex
```

merge(): Tipos de JOIN

Tipos de JOIN (como SQL):

| Tipo | Descrição |
|-------|--|
| inner | Apenas registros com chave em ambos |
| left | Todos de left, matching de right |
| right | Todos de right, matching de left |
| outer | Todos de ambos (union) |
| cross | Produto cartesiano (todas combinações) |

Padrão:

- ▶ Por padrão, `merge()` usa `how= 'inner'`
- ▶ Apenas linhas com chave comum em ambos os DataFrames

merge(): Inner Join

</> Python

```
1 passageiros = pd.DataFrame({ # DataFrame de passageiros
2     'PassengerId': [1, 2, 3, 4],
3     'Name': ['John', 'Anna', 'Peter', 'Linda']
4 })
5 cabines = pd.DataFrame({ # DataFrame de cabines
6     'PassengerId': [1, 2, 5, 6],
7     'Cabin': ['A1', 'B2', 'C3', 'D4']
8 })
9 # Inner join (apenas IDs comuns: 1, 2)
10 resultado = pd.merge(passageiros, cabines,
11                       on='PassengerId', how='inner')
12 print(resultado)
13 #      PassengerId  Name Cabin
14 # 0                1  John   A1
15 # 1                2  Anna   B2
```

merge(): Left Join

</> Python

```
1 passageiros = pd.DataFrame({
2     'PassengerId': [1, 2, 3, 4],
3     'Name': ['John', 'Anna', 'Peter', 'Linda']
4 })
5
6 cabines = pd.DataFrame({
7     'PassengerId': [1, 2, 5, 6],
8     'Cabin': ['A1', 'B2', 'C3', 'D4']
9 })
```

merge(): Left Join (cont.)

</> Python

```
1 # Left join (todos de passageiros)
2 resultado = pd.merge(passageiros, cabines,
3                       on='PassengerId', how='left')
4 print(resultado)
5 #      PassengerId  Name Cabin
6 # 0              1  John   A1
7 # 1              2  Anna   B2
8 # 2              3  Peter  NaN
9 # 3              4  Linda  NaN
10
```


merge(): Right e Outer Join

</> Python

```
1 # Right join (todos de cabines)
2 resultado = pd.merge(passageiros, cabines,
3                       on='PassengerId', how='right')
4 # Resultado: PassengerId 1, 2, 5, 6
5 # Outer join (todos de ambos)
6 resultado = pd.merge(passageiros, cabines,
7                       on='PassengerId', how='outer')
8 print(resultado)
9 #      PassengerId  Name Cabin
10 # 0              1   John   A1
11 # 1              2   Anna   B2
12 # 2              3  Peter  NaN
13 # 3              4  Linda  NaN
14 # 4              5    NaN   C3
15 # 5              6    NaN   D4
```

merge(): Múltiplas Chaves

Python

```
1 df1 = pd.DataFrame({
2     'ID': [1, 2, 3],
3     'Date': ['2024-01-01', '2024-01-01', '2024-01-02'],
4     'Value1': [10, 20, 30]
5 })
6 df2 = pd.DataFrame({
7     'ID': [1, 2, 3],
8     'Date': ['2024-01-01', '2024-01-02', '2024-01-01'],
9     'Value2': [100, 200, 300]
10 })
11 # Merge por multiplas colunas
12 resultado = pd.merge(df1, df2, on=['ID', 'Date'])
13 print(resultado)
14 # Apenas linhas com ID E Date iguais
```

merge(): Colunas com Nomes Diferentes

</> Python

```
1 passageiros = pd.DataFrame({  
2     'PassengerId': [1, 2, 3],  
3     'Name': ['John', 'Anna', 'Peter']  
4 })  
5 tickets = pd.DataFrame({  
6     'TicketId': [1, 2, 4],  
7     'Price': [50, 100, 150]  
8 })  
9
```

merge(): Colunas com Nomes Diferentes (cont.)

</> Python

```
1 # Especificar colunas diferentes
2 resultado = pd.merge(passageiros, tickets,
3                       left_on='PassengerId',
4                       right_on='TicketId',
5                       how='left')
6 print(resultado)
7 #      PassengerId   Name  TicketId  Price
8 # 0              1   John         1.0   50.0
9 # 1              2   Anna         2.0  100.0
10 # 2              3  Peter         NaN    NaN
```

merge(): Sufixos para Colunas Duplicadas

</> Python

```
1 df1 = pd.DataFrame({
2     'ID': [1, 2, 3],
3     'Value': [10, 20, 30],
4     'Status': ['A', 'B', 'C']
5 })
6
7 df2 = pd.DataFrame({
8     'ID': [1, 2, 3],
9     'Value': [100, 200, 300],
10    'Status': ['X', 'Y', 'Z']
11 })
12
```

merge(): Sufixos para Colunas Duplicadas (cont.)

</> Python

```
1 # Colunas 'Value' e 'Status' existem em ambos
2 resultado = pd.merge(df1, df2, on='ID',
3                       suffixes=('_left', '_right'))
4 print(resultado.columns)
5 # ['ID', 'Value_left', 'Status_left', 'Value_right', 'Status_right']
6
```

merge(): Indicador de Origem

</> Python

```
1 df1 = pd.DataFrame({'ID': [1, 2, 3], 'Val': [10, 20, 30]})
2 df2 = pd.DataFrame({'ID': [2, 3, 4], 'Val': [200, 300, 400]})
3 # Adicionar indicador de origem
4 resultado = pd.merge(df1, df2, on='ID', how='outer',
5                       indicator=True, suffixes=('_1', '_2'))
6 print(resultado)
7 #      ID  Val_1  Val_2    _merge
8 # 0     1   10.0   NaN  left_only
9 # 1     2   20.0  200.0     both
10 # 2     3   30.0  300.0     both
11 # 3     4    NaN  400.0  right_only
12 # Renomear indicador
13 resultado = pd.merge(df1, df2, on='ID', how='outer',
14                       indicator='source')
15 # Val_1 => Val_x, Val_2 => Val_y, _merge => source
```

Merge Validation: Evitar Erros

Problema comum:

- ▶ Merge pode duplicar linhas sem avisar
- ▶ Chaves duplicadas causam explosão de dados
- ▶ Difícil detectar em datasets grandes

Solução: Parâmetro validate:

- ▶ **'1:1'** - Um para um (ambas únicas)
- ▶ **'1:m'** - Um para muitos (left única, right pode repetir)
- ▶ **'m:1'** - Muitos para um (left pode repetir, right única)
- ▶ **'m:m'** - Muitos para muitos (ambas podem repetir)

Atenção

Sempre use validate em produção - previne bugs silenciosos!

validate: Exemplo de Erro

</> Python

```
1 # Clientes (deveria ter IDs únicos, mas tem duplicata!)
2 clientes = pd.DataFrame({
3     'id': [1, 2, 2, 3], # ID 2 duplicado!
4     'nome': ['Ana', 'Bruno', 'Bruno Jr', 'Carlos']
5 })
6 # Pedidos
7 pedidos = pd.DataFrame({
8     'id': [1, 2, 3],
9     'valor': [100, 200, 150]
10 })
11 # SEM validação: merge aceita, mas duplica linhas!
12 resultado_ruim = pd.merge(clientes, pedidos, on='id')
13 print(len(resultado_ruim)) # 4 linhas (deveria ser 3!)
```

validate: Detectar Erro

</> Python

```
1 clientes = pd.DataFrame({
2     'id': [1, 2, 2, 3],
3     'nome': ['Ana', 'Bruno', 'Bruno Jr', 'Carlos']
4 })
5 pedidos = pd.DataFrame({
6     'id': [1, 2, 3],
7     'valor': [100, 200, 150]
8 })
9 # COM validação: detecta problema!
10 try:
11     resultado = pd.merge(clientes, pedidos, on='id',
12                           validate='1:1')
13 except pd.errors.MergeError as e:
14     print(f"ERRO DETECTADO: {e}")
15     # MergeError: Merge keys are not unique in left dataset
```

validate: Casos de Uso

</> Python

```
1 # 1:1 - Enriquecer cadastro (ambas únicas)
2 clientes = pd.DataFrame({'id': [1, 2, 3], 'nome': ['A', 'B', 'C']})
3 enderecos = pd.DataFrame({'id': [1, 2, 3], 'cidade': ['SP', 'RJ', 'MG']})
4 resultado = pd.merge(clientes, enderecos, on='id', validate='1:1')
5
6 # 1:m - Um cliente, múltiplos pedidos
7 clientes = pd.DataFrame({'id': [1, 2], 'nome': ['Ana', 'Bruno']})
8 pedidos = pd.DataFrame({'id': [1, 1, 2], 'valor': [100, 150, 200]})
9 resultado = pd.merge(clientes, pedidos, on='id', validate='1:m')
10
11 # m:1 - Múltiplos produtos, uma categoria cada
12 produtos = pd.DataFrame({'cat_id': [1, 1, 2], 'prod': ['A', 'B', 'C']})
13 categorias = pd.DataFrame({'cat_id': [1, 2], 'cat': ['Eletro', 'Livros']})
14 resultado = pd.merge(produtos, categorias, on='cat_id', validate='m:1')
15
```

validate: Best Practices

</> Python

```
1 # SEMPRE especifique validate quando souber a relação
2 df1 = pd.read_csv('customers.csv')
3 df2 = pd.read_csv('orders.csv')
4 # BOM: Valida expectativa
5 resultado = pd.merge(
6     df1, df2,
7     on='customer_id',
8     how='left',
9     validate='1:m', # Um cliente, múltiplos pedidos
10    indicator=True  # Também adicionar indicador
11 )
12 # Verificar resultado
13 print(f"Total de clientes: {df1.shape[0]}")
14 print(f"Total após merge: {resultado.shape[0]}")
15 print(f"Distribuição: {resultado['_merge'].value_counts()}")
```

validate: Matriz de Decisão

| Relação | validate | Exemplo |
|----------------------|----------|--------------------------|
| Cliente ↔ Endereço | '1:1' | Cadastros complementares |
| Categoria → Produtos | '1:m' | Um para muitos |
| Pedidos → Categoria | 'm:1' | Muitos para um |
| Produtos ↔ Tags | 'm:m' | Muitos para muitos |

Quando não usar:

- ▶ Merge exploratório (ainda descobrindo relações)
- ▶ Dados sujos que precisam limpeza primeiro
- ▶ Relação é realmente m:m e você espera isso

💡 Nota Importante

Em produção: **sempre** use validate + indicator

merge_asof(): Merge por Proximidade

O que é merge_asof()?

- ▶ "Merge by nearest key" - chave mais próxima
- ▶ Útil para séries temporais
- ▶ Junta valores "aproximados" em vez de exatos
- ▶ Similar ao SQL ASOF JOIN

Casos de uso:

- ▶ Juntar preços de ações com eventos de notícias
- ▶ Correlacionar logs com timestamps diferentes
- ▶ Associar medições com tempo mais próximo
- ▶ Dados financeiros (quotes + trades)

💡 Nota Importante

Requer dados ordenados pela coluna de merge!

merge_asof(): Exemplo Básico

</> Python

```
1 # Cotações de ações (atualizadas a cada minuto)
2 quotes = pd.DataFrame({
3     'time': pd.to_datetime(['09:00', '09:01', '09:02', '09:03'], format='%
4     H:%M'),
5     'price': [100, 102, 101, 103]
6 })
7
8 # Trades (executadas em horários variados)
9 trades = pd.DataFrame({
10     'time': pd.to_datetime(['09:00:30', '09:01:45', '09:02:10'], format='%
11     H:%M:%S'),
12     'volume': [50, 100, 75]
13 })
```

merge_asof(): Exemplo Básico

</> Python

```
1 # Merge asof: pega preço mais recente para cada trade
2 resultado = pd.merge_asof(
3     trades,      # left (ordenado por time)
4     quotes,      # right (ordenado por time)
5     on='time',   # coluna temporal
6     direction='backward' # backward, forward, ou nearest
7 )
```


merge_asof(): Exemplo Básico(cont.)

</> Python

```
1 print(resultado)
2 #
3 # 0 2024-01-01 09:00:30      50      100 # Usou cotação 09:00
4 # 1 2024-01-01 09:01:45     100      102 # Usou cotação 09:01
5 # 2 2024-01-01 09:02:10      75      101 # Usou cotação 09:02
6
7 # direction='backward': pega valor anterior mais próximo
8 # direction='forward': pega valor posterior mais próximo
9 # direction='nearest': pega o mais próximo (qualquer direção)
10 # Com tolerância máxima
11 resultado_tol = pd.merge_asof(
12     trades, quotes,
13     on='time',
14     tolerance=pd.Timedelta('1min') # No máximo 1min de diferença
15 )
```

merge_asof(): Com Grupos

</> Python

```
1 # Cotações de múltiplas ações
2 quotes = pd.DataFrame({
3     'time': pd.to_datetime(['09:00', '09:01', '09:00', '09:01'], format='%H:%M'),
4     'ticker': ['AAPL', 'AAPL', 'GOOGL', 'GOOGL'],
5     'price': [150, 152, 2800, 2805]
6 })
7
8 # Trades
9 trades = pd.DataFrame({
10     'time': pd.to_datetime(['09:00:30', '09:00:45'], format='%H:%M:%S'),
11     'ticker': ['AAPL', 'GOOGL'],
12     'volume': [100, 50]
13 })
```

merge_asof(): Com Grupos (cont.)

</> Python

```
1 # Merge asof por ticker
2 resultado = pd.merge_asof(
3     trades.sort_values('time'),
4     quotes.sort_values('time'),
5     on='time',
6     by='ticker' # Agrupar por ticker
7 )
8 print(resultado)
9 #
10 #   time ticker volume price
11 # 0 1900-01-01 09:00:30 AAPL    100    150
12 # 1 1900-01-01 09:00:45 GOOGL    50   2800
```

explode(): Expandir Listas em Linhas

O que é explode()?

- ▶ Transforma valores tipo lista em linhas separadas
- ▶ "Explode" uma linha em múltiplas
- ▶ Útil para dados aninhados
- ▶ Repete valores das outras colunas

Casos de uso:

- ▶ Dados de e-commerce (produtos por pedido)
- ▶ Tags de artigos/posts
- ▶ Dados JSON aninhados
- ▶ Relacionamentos muitos-para-muitos

💡 Nota Importante

Oposto de `groupby().agg(list)`

explode(): Exemplo

</> Python

```
1 # Pedidos com múltiplos produtos
2 pedidos = pd.DataFrame({
3     'pedido_id': [1, 2, 3],
4     'cliente': ['Ana', 'Bruno', 'Carlos'],
5     'produtos': [
6         ['Mouse', 'Teclado'],
7         ['Monitor'],
8         ['Webcam', 'Headset', 'Mouse']
9     ]
10 })
11 print(pedidos)
12 #      pedido_id  cliente      produtos
13 # 0           1      Ana    [Mouse, Teclado]
14 # 1           2     Bruno      [Monitor]
15 # 2           3    Carlos [Webcam, Headset, Mouse]
```

explode(): Exemplo (cont.)

</> Python

```
1 # Explodir coluna de produtos
2 resultado = pedidos.explode('produtos')
3 print(resultado)
4 #      pedido_id  cliente produtos
5 # 0             1      Ana    Mouse
6 # 0             1      Ana  Teclado
7 # 1             2     Bruno  Monitor
8 # 2             3    Carlos   Webcam
9 # 2             3    Carlos  Headset
10 # 2             3    Carlos    Mouse
11 # Note: índice repete para linhas do mesmo pedido
12 # Resetar índice se necessário
13 resultado_limpo = resultado.reset_index(drop=True)
```

explode(): Múltiplas Colunas

</> Python

```
1 # Explodir múltiplas colunas simultaneamente
2 pedidos_completo = pd.DataFrame({
3     'pedido_id': [1, 2],
4     'produtos': [['Mouse', 'Teclado'], ['Monitor']],
5     'precos': [[45, 120], [890]]
6 })
7 resultado = pedidos_completo.explode(['produtos', 'precos'])
8 resultado_limpo = resultado.reset_index(drop=True)
9 print(resultado_limpo)
10 #      pedido_id  produtos  precos
11 # 0             1    Mouse     45
12 # 1             1  Teclado    120
13 # 2             2  Monitor    890
```

explode(): Múltiplas Colunas

</> Python

```
1 # Agora podemos calcular total por pedido
2 total = resultado.groupby('pedido_id')['precos'].sum()
3 print(total)
4 # pedido_id
5 # 1      165
6 # 2      890
7 Name: precos, dtype: object
```


explode(): Pipeline Completo

</> Python

```
1 # Caso real: análise de tags de artigos
2 artigos = pd.DataFrame({
3     'titulo': ['Post 1', 'Post 2', 'Post 3'],
4     'autor': ['Ana', 'Bruno', 'Ana'],
5     'tags': [
6         ['python', 'data'],
7         ['sql', 'data', 'database'],
8         ['python', 'ml']
9     ]
10 })
```

explode(): Pipeline Completo (cont.)

</> Python

```
1 # Pipeline: explodir e contar
2 analise = (
3     artigos
4     .explode('tags')
5     .groupby('tags')
6     .agg(
7         total_artigos=('titulo', 'count'),
8         autores=('autor', lambda x: ', '.join(x.unique()))
9     )
10    .sort_values('total_artigos', ascending=False)
11 )
12 print(analise)
```

join(): Merge pelo Índice

</> Python

```
1 df1 = pd.DataFrame({
2     'A': [1, 2, 3]
3 }, index=['x', 'y', 'z'])
4 df2 = pd.DataFrame({
5     'B': [4, 5, 6]
6 }, index=['x', 'y', 'w'])
7 resultado = df1.join(df2) # Join pelo indice (left por padrao)
8 print(resultado)
9 #      A      B
10 # x    1    4.0
11 # y    2    5.0
12 # z    3    NaN
13 # Equivalente a merge com left_index e right_index
14 resultado = pd.merge(df1, df2, left_index=True,
15                       right_index=True, how='left')
```

Exemplo Prático: Enriquecer Dados

</> Python

```
1 # Dataset principal
2 df = pd.read_csv('titanic.csv')
3 # Informacoes adicionais de portos
4 portos = pd.DataFrame({
5     'Embarked': ['S', 'C', 'Q'],
6     'PortoNome': ['Southampton', 'Cherbourg', 'Queenstown'],
7     'Pais': ['Inglaterra', 'Franca', 'Irlanda']
8 })
9 # Adicionar informacoes dos portos
10 df_enriquecido = pd.merge(df, portos, on='Embarked',
11                             how='left')
12
13 print(df_enriquecido[['Name', 'Embarked', 'PortoNome']].head())
```

Exemplo Prático: Combinar Dados Temporais

Python

```
1 # Dados de diferentes periodos
2 jan = pd.DataFrame({
3     'ID': [1, 2, 3],
4     'VendasJan': [100, 200, 150]
5 })
6 fev = pd.DataFrame({
7     'ID': [1, 2, 3],
8     'VendasFev': [120, 180, 160]
9 })
10 # Combinar horizontalmente
11 vendas = pd.merge(jan, fev, on='ID')
12 # Calcular total
13 vendas['Total'] = vendas['VendasJan'] + vendas['VendasFev']
14 print(vendas)
```

Bloco 4

Reshape e Pivot

Reshape: Reorganizar Dados

Por que reestruturar dados?

- ▶ Formatos diferentes para análises diferentes
- ▶ Preparar dados para visualização
- ▶ Converter entre formato wide e long
- ▶ Criar tabelas de contingência
- ▶ Facilitar agregações específicas

Principais operações:

- ▶ **pivot_table()**: Wide format (tabela de resumo)
- ▶ **melt()**: Long format (unpivot)
- ▶ **stack()/unstack()**: Reorganizar índices
- ▶ **crosstab()**: Tabulação cruzada

Wide vs Long Format

Wide Format (formato largo):

- ▶ Uma linha por entidade
- ▶ Múltiplas colunas para valores diferentes
- ▶ Boa para leitura humana
- ▶ Exemplo: Notas de alunos (uma coluna por prova)

Long Format (formato longo):

- ▶ Múltiplas linhas por entidade
- ▶ Colunas: ID, variável, valor
- ▶ Melhor para análise e visualização
- ▶ Exemplo: Notas com colunas [Aluno, Prova, Nota]

Nota Importante

Muitas funções do Pandas preferem long format

pivot_table(): Criar Tabela de Resumo

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Pivot table basica: sobrevivencia por classe e sexo
3 pivot = pd.pivot_table(
4     df,
5     values='Survived',
6     index='Pclass',
7     columns='Sex',
8     aggfunc='mean'
9 )
10 print(pivot)
11 # Sex      female      male
12 # Pclass
13 # 1      0.968085  0.368852
14 # 2      0.921053  0.157407
15 # 3      0.500000  0.135447
```

pivot_table(): Múltiplas Agregações

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 pivot = pd.pivot_table( # Múltiplas funções de agregação
3     df,
4     values='Fare',
5     index='Pclass',
6     columns='Sex',
7     aggfunc=['mean', 'median', 'count']
8 )
9 print(pivot)
10 # MultiIndex em colunas
11 # mean          median          count
12 # Sex  female male  female male  female male
13 # ...
14 # Achatar colunas
15 pivot.columns = ['_'.join(col) for col in pivot.columns]
```

pivot_table(): Múltiplos Valores

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 pivot = pd.pivot_table(    # Agregar multiples columnas
3     df,
4     values=['Survived', 'Age'],
5     index='Pclass',
6     columns='Sex',
7     aggfunc='mean'
8 )
9 print(pivot)
```

| # | Age | Survived | | |
|----------|--------|----------|--------|-------|
| # Sex | female | male | female | male |
| # Pclass | | | | |
| # 1 | 34.61 | 41.28 | 0.968 | 0.369 |
| # 2 | 28.72 | 30.74 | 0.921 | 0.157 |
| # 3 | 21.75 | 26.51 | 0.500 | 0.135 |

pivot_table(): Totais (Margins)

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Adicionar linha/coluna de totais
4 pivot = pd.pivot_table(
5     df,
6     values='Survived',
7     index='Pclass',
8     columns='Sex',
9     aggfunc='mean',
10    margins=True,
11    margins_name='Total'
12 )
13
```

pivot_table(): Totais (Margins) (cont.)

</> Python

```
1 print(pivot)
2 # Sex      female      male      Total
3 # Pclass
4 # 1      0.968085    0.368852    0.629630
5 # 2      0.921053    0.157407    0.472826
6 # 3      0.500000    0.135447    0.242363
7 # Total   0.742038    0.188908    0.383838
8
```

pivot_table(): Preencher Valores Ausentes

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 pivot = pd.pivot_table( # Algumas combinacoes podem nao existir
3     df,
4     values='Fare',
5     index='Pclass',
6     columns='Embarked',
7     aggfunc='mean',
8     fill_value=0 # Preencher NaN com 0
9 )
10 print(pivot)
11 # Embarked          C          Q          S
12 # Pclass
13 # 1          104.718529   90.000000   70.364862
14 # 2          25.358335   12.350000   20.327439
15 # 3          11.214083   11.183393   14.644083
```

pivot_table(): Preencher Valores Ausentes (cont.)

</> Python

```
1 # Alternativa: dropna
2 pivot = pd.pivot_table(
3     df,
4     values='Fare',
5     index='Pclass',
6     columns='Embarked',
7     aggfunc='mean',
8     dropna=False
9 )
10 print(pivot)
11 # Embarked          C          Q          S    NaN
12 # Pclass
13 # 1      104.718529   90.000000   70.364862   80.0
14 # 2      25.358335   12.350000   20.327439    NaN
15 # 3      11.214083   11.183393   14.644083    NaN
```

crosstab(): Tabulação Cruzada

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Tabela de frequencias
3 tabela = pd.crosstab(df['Pclass'], df['Sex'])
4 print(tabela)
5 # Sex      female  male
6 # Pclass
7 # 1           94   122
8 # 2           76   108
9 # 3          144   347
10 # Com percentuais
11 tabela_pct = pd.crosstab(df['Pclass'], df['Sex'],
12                          normalize='index')
13 print(tabela_pct)
14 # Proporcoes por linha (soma = 1.0 em cada linha)
```


crosstab(): Com Valores

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Crosstab com valores agregados
3 tabela = pd.crosstab(
4     df['Pclass'],
5     df['Sex'],
6     values=df['Fare'],
7     aggfunc='mean'
8 )
9 print(tabela)
10 # Tarifa media por classe e sexo, single index
```

crosstab(): Com Valores (cont.)

</> Python

```
1 # Com multiplas funcoes
2 tabela = pd.crosstab(
3     df['Pclass'],
4     df['Sex'],
5     values=df['Fare'],
6     aggfunc=['mean', 'median', 'count'] #MultiIndex
7 )
8 print(tabela)
```

| # | | mean | | median | | count | |
|---|--------|------------|-----------|----------|---------|--------|------|
| # | Sex | female | male | female | male | female | male |
| # | Pclass | | | | | | |
| # | 1 | 106.125798 | 67.226127 | 82.66455 | 41.2625 | 94 | 122 |
| # | 2 | 21.970121 | 19.741782 | 22.00000 | 13.0000 | 76 | 108 |
| # | 3 | 16.118810 | 12.661633 | 12.47500 | 7.9250 | 144 | 347 |

crosstab(): Múltiplas Variáveis

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Crosstab com multiplas variaveis em linhas/colunas
4 tabela = pd.crosstab(
5     [df['Pclass'], df['Sex']],
6     [df['Survived'], df['Embarked']]
7 )
8 print(tabela)
9 # MultiIndex em linhas e colunas
10
```

crosstab(): Múltiplas Variáveis (cont.)

</> Python

```
1 # Com totais
2 tabela = pd.crosstab(
3     df['Pclass'],
4     df['Sex'],
5     margins=True
6 )
7 print(tabela)
8 # Linha e coluna 'All' com totais
9
```

melt(): Wide to Long (Unpivot)

</> Python

```
1 # Dados wide
2 df_wide = pd.DataFrame({
3     'Name': ['Ana', 'Bruno', 'Carlos'],
4     'Math': [85, 90, 78],
5     'Physics': [88, 92, 80],
6     'Chemistry': [90, 85, 82]
7 })
8 print(df_wide)
9 #      Name  Math  Physics  Chemistry
10 # 0     Ana    85      88         90
11 # 1    Bruno    90      92         85
12 # 2   Carlos    78      80         82
13
14
```

melt(): Wide to Long (Unpivot) (cont.)

</> Python

```
1 # Converter para long
2 df_long = pd.melt(
3     df_wide,
4     id_vars=['Name'],
5     var_name='Subject',
6     value_name='Grade'
7 )
8 print(df_long)
9 #      Name      Subject  Grade
10 # 0      Ana      Math      85
11 # 1     Bruno      Math      90
12 # 2    Carlos      Math      78
13 # 3      Ana    Physics      88
14 # ...
```

melt(): Selecionar Colunas

</> Python

```
1 df_wide = pd.DataFrame({
2     'ID': [1, 2, 3],
3     'Name': ['Ana', 'Bruno', 'Carlos'],
4     'Q1': [100, 150, 120],
5     'Q2': [110, 160, 130],
6     'Q3': [105, 155, 125],
7     'Q4': [115, 165, 135]
8 })
9 df_long = pd.melt( # Melt apenas colunas de trimestres
10     df_wide,
11     id_vars=['ID', 'Name'],
12     value_vars=['Q1', 'Q2', 'Q3', 'Q4'],
13     var_name='Quarter',
14     value_name='Sales'
15 )
```

stack() e unstack(): Reorganizar Níveis

</> Python

```
1 # DataFrame com MultiIndex
2 df = pd.DataFrame({
3     'A': [1, 2, 3, 4],
4     'B': [5, 6, 7, 8]
5 }, index=pd.MultiIndex.from_tuples([
6     ('X', 'a'), ('X', 'b'), ('Y', 'a'), ('Y', 'b')
7 ]))
8
9 # Unstack: mover nível do índice para colunas
10 df_unstacked = df.unstack()
11 print(df_unstacked)
```

| # | | A | | B | |
|----|-----|---|---|---|---|
| # | | a | b | a | b |
| 14 | # X | 1 | 2 | 5 | 6 |
| 15 | # Y | 3 | 4 | 7 | 8 |

stack() e unstack(): Reorganizar Níveis (cont.)

</> Python

```
1 # Stack: mover nível das colunas para índice
2 df_stacked = df_unstacked.stack()
3 print(df_stacked)
4 #      A  B
5 # X a  1  5
6 #   b  2  6
7 # Y a  3  7
8 #   b  4  8
```

stack() e unstack(): Exemplo Prático

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # Criar pivot simples
4 pivot = df.pivot_table(
5     values='Survived',
6     index='Pclass',
7     columns='Sex',
8     aggfunc='mean'
9 )
10
```

stack() e unstack(): Exemplo Prático (cont.)

</> Python

```
1 # Unstack: transformar colunas em linhas
2 unstacked = pivot.unstack()
3 print(unstacked)
4 # Sex      Pclass
5 # female  1      0.968085
6 #         2      0.921053
7 #         3      0.500000
8 # male    1      0.368852
9 #         2      0.157407
10 #        3      0.135447
11
```

Tratamento de Duplicatas

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Verificar duplicatas
3 duplicatas = df.duplicated()
4 print(f"Linhas duplicadas: {duplicatas.sum()}")
5 # Duplicatas em colunas especificas
6 dup_nome = df.duplicated(subset=['Name'])
7 print(f"Nomes duplicados: {dup_nome.sum()}")
8 # Ver as duplicatas
9 df_duplicados = df[df.duplicated(keep=False)]
10 # Manter apenas primeira ocorrencia
11 df_limpo = df.drop_duplicates()
12 # Manter ultima ocorrencia
13 df_limpo = df.drop_duplicates(keep='last')
```

drop_duplicates(): Colunas Específicas

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Remover duplicatas baseado em colunas especificas
3 # Manter primeiro passageiro de cada combinacao classe+sexo
4 df_unico = df.drop_duplicates(subset=['Pclass', 'Sex'])
5 print(f"Original: {len(df)}, Unico: {len(df_unico)}")
6 # Original: 891, Unico: 6
7 # Verificar combinacoes unicas
8 combinacoes = df[['Pclass', 'Sex']].drop_duplicates()
9 print(combinacoes)
10 # Indice, Pclass, Sex - indices nao-consecutivos
11 # Contar ocorrencias de cada combinacao
12 contagem = df.groupby(['Pclass', 'Sex']).size()
13 print(contagem)
14 # Numero de ocorrencias de cada combinacao
```

Exemplo Completo: Pipeline de Transformação

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # 1. Remover duplicatas
3 df = df.drop_duplicates()
4 # 2. Criar novas features
5 df['FamilySize'] = df['SibSp'] + df['Parch'] + 1
6 df['IsAlone'] = (df['FamilySize'] == 1).astype(int)
7 df['Titulo'] = df['Name'].str.extract(r'([^.]+\.)\.')
8 # 3. Agrupar titulos raros
9 titulos_comuns = ['Mr', 'Miss', 'Mrs', 'Master']
10 df['TituloGroup'] = df['Titulo'].apply(
11     lambda x: x if x in titulos_comuns else 'Other'
12 )
```

Exemplo Completo: Pipeline (cont.)

</> Python

```
1 # 4. Preencher valores faltantes por grupo
2 df['Age'] = df.groupby('TituloGroup')['Age'].transform(
3     lambda x: x.fillna(x.median())
4 )
5 # 5. Criar tabela de resumo
6 resumo = pd.pivot_table(
7     df,
8     values=['Survived', 'Age', 'Fare'],
9     index='Pclass',
10    columns='Sex',
11    aggfunc={'Survived': 'mean',
12             'Age': 'median',
13             'Fare': 'mean'}
14 )
```

Exemplo Completo: Pipeline (cont.)

</> Python

```
1 print(resumo)
2 #
3 # Sex      female  male      Fare      Survived
4 # Pclass
5 # 1      35.0    36.5    106.125798    67.226127    0.968085    0.368852
6 # 2      28.0    30.0     21.970121    19.741782    0.921053    0.157407
7 # 3      21.0    30.0     16.118810    12.661633    0.500000    0.135447
```


Exemplo Prático: Análise de Família

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Criar features de familia
3 df['FamilySize'] = df['SibSp'] + df['Parch'] + 1
4 df['FamilyType'] = pd.cut(
5     df['FamilySize'],
6     bins=[0, 1, 4, 20],
7     labels=['Alone', 'Small', 'Large']
8 )
9 # Analise por tipo de familia
10 analise = df.groupby('FamilyType', observed=False).agg({
11     'Survived': ['count', 'sum', 'mean'],
12     'Age': 'mean',
13     'Fare': 'mean'
14 })
15 print(analise)
```

Exemplo Prático: Enriquecer com Estatísticas

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Adicionar estatísticas do grupo como novas colunas
3 df['AgeMeanClass'] = df.groupby('Pclass')['Age'].transform('mean')
4 df['AgeStdClass'] = df.groupby('Pclass')['Age'].transform('std')
5 # Z-score normalizado por classe
6 df['AgeZScore'] = (df['Age'] - df['AgeMeanClass']) / \
7                   df['AgeStdClass']
8 # Ranking de tarifa dentro da classe
9 df['FareRank'] = df.groupby('Pclass')['Fare'].rank(
10     ascending=False
11 )
12 print(df[['Name', 'Pclass', 'Age', 'AgeMeanClass', 'AgeZScore']].head())
13 # Adiciona colunas com estatísticas (repetidas) para cada linha
```

Exemplo Prático: Relatório Multi-nível

Python

```
1 df = pd.read_csv('titanic.csv')
2 # Criar relatório complexo
3 relatorio = df.groupby(['Pclass', 'Sex', 'Embarked']).agg({
4     'PassengerId': 'count',
5     'Survived': ['sum', 'mean'],
6     'Age': ['mean', 'median'],
7     'Fare': ['mean', 'max']
8 })
9 # Achatar MultiIndex de colunas
10 relatorio.columns = ['_'.join(col) for col in relatorio.columns]
11 # Resetar índice para facilitar visualização
12 relatorio = relatorio.reset_index()
```

Exemplo Prático: Relatório Multi-nível (cont.)

</> Python

```
1 print(relatorio.columns)
2 # Index(['Pclass', 'Sex', 'Embarked', 'PassengerId_count', 'Survived_sum',
3 #       'Survived_mean', 'Age_mean', 'Age_median', 'Fare_mean', 'Fare_max
4 #       '],
5 #       dtype='object')
```

Method Chaining: Código Limpo e Legível

O que é Method Chaining?

- ▶ Encadear operações em sequência
- ▶ Cada método retorna DataFrame/Series
- ▶ Uma operação por linha
- ▶ Estilo "fluent interface"

Vantagens:

- ▶ Código mais legível (lê-se como pipeline)
- ▶ Menos variáveis intermediárias
- ▶ Fácil adicionar/remover etapas
- ▶ Documentação clara do fluxo
- ▶ Evita modificações acidentais

💡 Nota Importante

Estilo recomendado para código de produção

Sem vs Com Method Chaining

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # SEM method chaining (múltiplas variáveis)
4 df_limpo = df.dropna(subset=['Age'])
5 df_filtrado = df_limpo[df_limpo['Pclass'] == 1]
6 df_com_features = df_filtrado.copy()
7 df_com_features['FamilySize'] = df_com_features['SibSp'] + df_com_features
   ['Parch'] + 1
8 df_ordenado = df_com_features.sort_values('Fare', ascending=False)
9 resultado = df_ordenado.head(10)
```

Sem vs Com Method Chaining

</> Python

```
1 df = pd.read_csv('titanic.csv')
2
3 # COM method chaining (pipeline claro)
4 resultado = (
5     df
6     .dropna(subset=['Age'])
7     .query('Pclass == 1')
8     .assign(FamilySize=lambda x: x['SibSp'] + x['Parch'] + 1)
9     .sort_values('Fare', ascending=False)
10    .head(10)
11 )
```

assign(): Criar Colunas no Pipeline

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # assign() retorna cópia com novas colunas
3 resultado = (
4     df
5     .assign(
6         FamilySize=lambda x: x['SibSp'] + x['Parch'] + 1, # n colunas
7         IsAlone=lambda x: x['FamilySize'] == 1,
8         FamilyType=lambda x: pd.cut( # Usar coluna recém-criada
9             x['FamilySize'],
10            bins=[0, 1, 4, 20],
11            labels=['Alone', 'Small', 'Large']
12        )
13    )
14    .head()
15 )
```


assign(): Criar Colunas no Pipeline (cont.)

</> Python

```
1 print(resultado[['Name', 'FamilySize', 'IsAlone', 'FamilyType']])
2 #
3 # 0 Braund, Mr. Owen                2    False    Small
4 # 1 Cumings, Mrs. [...]            2    False    Small
5 # 2 Heikkinen, [...]               1     True    Alone
6 # 3 Futrelle, [...]                2    False    Small
7 # 4 Allen, Mr. [...]               1     True    Alone
```

pipe(): Funções Customizadas no Pipeline

</> Python

```
1 # Definir funções de transformação
2 def limpar_dados(df):
3     return df.dropna(subset=['Age', 'Embarked'])
4
5 def criar_features(df):
6     return df.assign(
7         FamilySize=lambda x: x['SibSp'] + x['Parch'] + 1,
8         AgeGroup=lambda x: pd.cut(x['Age'], bins=[0, 18, 60, 100],
9                                   labels=['Child', 'Adult', 'Senior'])
10    )
11
12 def filtrar_primeira_classe(df):
13     return df.query('Pclass == 1')
```

pipe(): Funções Customizadas no Pipeline (cont.)

</> Python

```
1 # Pipeline com pipe()
2 resultado = (
3     pd.read_csv('titanic.csv')
4     .pipe(limpar_dados)
5     .pipe(criar_features)
6     .pipe(filtrar_primeira_classe)
7     .head()
8 )
```

Pipeline Completo: ETL

</> Python

```
1 def processar_titanic(filepath):
2     return (
3         # EXTRACT
4         pd.read_csv(filepath)
5         # TRANSFORM - Limpeza
6         .dropna(subset=['Embarked'])
7         .drop_duplicates()
8         .assign(Age=lambda x: x.groupby('Pclass')['Age']
9                 .transform(lambda y: y.fillna(y.median()))))
10    # TRANSFORM - Features
11    .assign(
12        FamilySize=lambda x: x['SibSp'] + x['Parch'] + 1,
13        Title=lambda x: x['Name'].str.extract(r'([^.]+\.)')[0],
14        AgeGroup=lambda x: pd.cut(x['Age'], bins=[0,18,60,100],
15                                  labels=['Child', 'Adult', 'Senior'])
16    )
```

Pipeline Completo: ETL (cont.)

</> Python

```
1      # TRANSFORM - Tipos
2      .astype({
3          'Pclass': 'category',
4          'Sex': 'category',
5          'Embarked': 'category',
6          'AgeGroup': 'category'
7      })
8      # TRANSFORM - Ordenação
9      .sort_values(['Pclass', 'Fare'], ascending=[True, False])
10     .reset_index(drop=True)
11 )
12 # Executar pipeline
13 df_processado = processar_titanic('titanic.csv')
14 # Pipeline é reutilizável e testável!
```

Debugging em Pipelines

</> Python

```
1 # Técnica 1: Comentar etapas temporariamente
2 resultado = (
3     df
4     .dropna()
5     .query('Age > 18')
6     # .assign(NewCol=lambda x: x['A'] * 2) # Comentar para debug
7     .head()
8 )
```

Debugging em Pipelines

</> Python

```
1 # Técnica 2: Quebrar pipeline com variável intermediária
2 temp = (
3     df
4     .dropna()
5     .query('Age > 18')
6 )
7 print(temp.shape) # Verificar até aqui
8 resultado = temp.assign(NewCol=lambda x: x['A'] * 2)
```

Debugging em Pipelines

</> Python

```
1 # Técnica 3: Usar pipe() para debug
2 def debug_shape(df):
3     print(f"Shape atual: {df.shape}")
4     return df
5
6 resultado = df.pipe(debug_shape).dropna().pipe(debug_shape)
```


Method Chaining: Best Practices

Regras de ouro:

1. **Uma operação por linha**
 - ▶ Facilita leitura e debug
2. **Use parênteses externos**
 - ▶ Permite quebra de linha sem \
3. **Indentação consistente**
 - ▶ 4 espaços para cada nível
4. **assign() para novas colunas**
 - ▶ Em vez de `df['col'] = valor`
5. **query() para filtros simples**
 - ▶ Mais legível que boolean indexing
6. **pipe() para funções complexas**
 - ▶ Mantém pipeline fluido

Exemplo Real: Análise de Vendas

</> Python

```
1 # Pipeline completo de análise
2 analise_vendas = (
3     pd.read_csv('vendas.csv')
4     # Limpeza
5     .dropna(subset=['valor', 'data'])
6     .drop_duplicates(subset=['pedido_id'])
7     # Datas
8     .assign(data=lambda x: pd.to_datetime(x['data']))
9     .assign(
10         mes=lambda x: x['data'].dt.to_period('M'),
11         dia_semana=lambda x: x['data'].dt.day_name()
12     )
```

Exemplo Real: Análise de Vendas (cont.)

</> Python

```
1  .assign(  
2      receita=lambda x: x['quantidade'] * x['preco'],  
3      margem=lambda x: x['receita'] - x['custo']  
4  )  
5  # Agregação  
6  .groupby(['mes', 'categoria'])  
7  .agg(  
8      receita_total=('receita', 'sum'),  
9      margem_media=('margem', 'mean'),  
10     num_pedidos=('pedido_id', 'count')  
11 )  
12 .round(2)  
13 )
```

Performance: Otimização de Código Pandas

Hierarquia de performance (do mais rápido ao mais lento):

1. Operações vetorizadas (100x - 1000x mais rápido)

- ▶ `df['col'] * 2, df['A'] + df['B']`

2. Built-in Pandas functions (10x - 100x)

- ▶ `.groupby(), .merge(), .query()`

3. `apply()` com funções NumPy (5x - 10x)

- ▶ `df.apply(np.mean)`

4. `apply()` com Python puro (2x - 5x)

- ▶ `df['col'].apply(lambda x: x + 1)`

5. Loops explícitos (NUNCA USE!)

- ▶ `for i in range(len(df)): ...`

Performance: Exemplo Comparativo

</> Python

```
1 import time
2 df = pd.DataFrame({'valor': range(1000000)})
3 # LENTO: Loop explícito (NUNCA faça isso!)
4 start = time.time()
5 resultado = []
6 for i in range(len(df)):
7     resultado.append(df.iloc[i]['valor'] * 2)
8 df['loop'] = resultado
9 print(f"Loop: {time.time() - start:.2f}s") # ~5 segundos
```

Performance: Exemplo Comparativo (cont.)

</> Python

```
1 # MÉDIO: apply()
2 start = time.time()
3 df['apply'] = df['valor'].apply(lambda x: x * 2)
4 print(f"Apply: {time.time() - start:.2f}s") # ~0.5 segundos
5
6 # RÁPIDO: Vetorizado
7 start = time.time()
8 df['vetorizado'] = df['valor'] * 2
9 print(f"Vetorizado: {time.time() - start:.2f}s") # ~0.005s
```

Performance: Uso de Memória

</> Python

```
1 df = pd.read_csv('titanic.csv')
2 # Verificar uso de memória
3 print(df.memory_usage(deep=True))
4 # Index          128
5 # PassengerId    7128
6 # Name          52194  <- Strings pesam muito!
7 # Sex           5118
8 # ...
```

Performance: Uso de Memória (cont.)

</> Python

```
1 # ANTES: object dtype
2 print(f"Memória Sex: {df['Sex'].memory_usage(deep=True)} bytes")
3 # 52500 bytes
4
5 # DEPOIS: category dtype
6 df['Sex'] = df['Sex'].astype('category')
7 print(f"Memória Sex: {df['Sex'].memory_usage(deep=True)} bytes")
8 # 1019 bytes (50x menor!)
```


Performance: Dtypes Apropriados

Python

```
1 df = pd.read_csv('titanic.csv')
2 # Reduzir tamanho de inteiros/floats
3 print(df['PassengerId'].dtype) # int64
4 df['PassengerId'] = df['PassengerId'].astype('int16') # 0-32767
5 print(df['Age'].dtype) # float64
6 df['Age'] = df['Age'].astype('float32')
7 # Converter colunas repetitivas para category
8 for col in ['Pclass', 'Sex', 'Embarked', 'Survived']:
9     df[col] = df[col].astype('category')
10 # Comparar memória
11 print("\nEconomia de memória:")
12 original = pd.read_csv('titanic.csv')
13 print(f"Original: {original.memory_usage(deep=True).sum() / 1024:.1f} KB")
14 print(f"Otimizado: {df.memory_usage(deep=True).sum() / 1024:.1f} KB")
```

Performance: Filtrar Antes de Agrupar

</> Python

```
1 df = pd.read_csv('huge_dataset.csv') # 10 milhões de linhas
2 # LENTO: Agrupar tudo, depois filtrar
3 resultado_lento = (
4     df
5     .groupby('categoria')['valor'].sum()
6     .loc[lambda x: x > 1000] # Filtra depois
7 )
8 # RÁPIDO: Filtrar primeiro, depois agrupar
9 resultado_rapido = (
10     df
11     .query('valor > 10') # Reduz dados primeiro
12     .groupby('categoria')['valor'].sum()
13     .loc[lambda x: x > 1000]
14 )
15 # Reduz de 10M para 1M linhas ANTES de agrupar
16 # Economia: ~10x mais rápido
```

Performance: Profiling com %%timeit

</> Python

```
1 # No Jupyter/Colab: usar %%timeit para medir performance
2 %%timeit
3 df['resultado'] = df['A'] + df['B']
4 # 1.23 ms ± 45.2 µs per loop (mean ± std. dev. of 7 runs)
5 %%timeit
6 df['resultado'] = df.apply(lambda row: row['A'] + row['B'], axis=1)
7 # 456 ms ± 12.3 ms per loop
8 # -> 370x mais lento!
9 # Comparar abordagens
10 %timeit resultado1 = df['A'] * 2 # Vetorizado
11 %timeit resultado2 = df['A'].apply(lambda x: x * 2) # Apply
12 %timeit resultado3 = [x * 2 for x in df['A']] # List comp
13 # Regra: Sempre teste alternativas em dados grandes!
```

Performance: Checklist de Otimização

Antes de otimizar:

1. **Meça primeiro!** Use `%timeit` ou `time.time()`
2. Otimize apenas gargalos reais (regra 80/20)

Otimizações de alto impacto:

- ✓ Vetorizar em vez de `apply()` ou loops
- ✓ Usar dtypes apropriados (`category`, `int32`)
- ✓ Filtrar dados antes de operações pesadas
- ✓ Usar `query()` em vez de boolean indexing complexo
- ✓ Usar `inplace=True` com cuidado (economiza memória)

Trade-offs:

- ▶ Legibilidade vs Performance: Priorize legibilidade até medir gargalo
- ▶ Memória vs Velocidade: Categoricals economizam memória, podem ser mais lentos

Dicas Finais: Pandas em Produção

Boas práticas essenciais:

1. Validação de dados

- ▶ Sempre verifique shape após operações
- ▶ Use validate em merges
- ▶ Adicione indicator=True para rastrear

2. Tratamento de erros

- ▶ Try-except em leituras de arquivo
- ▶ Validar dtypes esperados
- ▶ Verificar valores faltantes críticos

3. Documentação

- ▶ Docstrings em funções de pipeline
- ▶ Comentar transformações complexas
- ▶ Registrar decisões de limpeza

4. Testes

- ▶ Testes unitários para funções de transformação
- ▶ Verificar invariantes (ex: sum antes/depois)

Dicas Finais: Debugging Comum

</> Python

```
1 # 1. KeyError: coluna não existe
2 try:
3     df['coluna_inexistente']
4 except KeyError:
5     print(f"Colunas disponíveis: {df.columns.tolist()}")
6
7 # 2. SettingWithCopyWarning
8 # Evite: df_subset['nova'] = valor
9 # Use: df.loc[mask, 'nova'] = valor
10
11 # 3. Merge duplicando linhas
12 resultado = pd.merge(df1, df2, on='id', validate='1:1')
```

Dicas Finais: Debugging Comum (cont.)

Python

```
1 # 4. GroupBy perdendo dados
2 antes = len(df)
3 depois = len(df.groupby('cat').size().sum())
4 assert antes == depois, "GroupBy perdeu linhas!"
5
6 # 5. Tipos incorretos após merge
7 df['id'] = df['id'].astype(int) # Forçar tipo correto
8
```

Recursos para Aprofundamento

Documentação oficial:

- ▶ User Guide: https://pandas.pydata.org/docs/user_guide/
- ▶ API Reference: Consulte sempre que precisar
- ▶ Performance: Enhancement proposals (PDEP)

Livros recomendados:

- ▶ "Python for Data Analysis- Wes McKinney (criador do Pandas)
- ▶ "Pandas 2.x Cookbook- Matt Harrison
- ▶ "Effective Pandas- Matt Harrison

Prática:

- ▶ 100 Pandas Puzzles (GitHub)
- ▶ Kaggle: Datasets + Notebooks da comunidade
- ▶ Real Python: Tutoriais práticos

Live Coding

Vamos praticar juntos:

1. Transformações com `apply()` e `map()`
2. Operações com strings (str accessor)
3. `GroupBy` com múltiplas agregações
4. Merge de DataFrames complementares
5. Criação de pivot tables
6. Pipeline completo de transformação

Revisão da Aula

Bloco 1 - Transformações:

- ▶ Adicionar/remover colunas
- ▶ `apply()`, `map()` para transformações customizadas
- ▶ String methods (str accessor)
- ▶ Conversão de tipos: `astype()`, `to_numeric()`
- ▶ Categorical para economia de memória

Bloco 2 - GroupBy:

- ▶ Split-apply-combine paradigm
- ▶ Agregações: `agg()`, `sum()`, `mean()`, `count()`
- ▶ `transform()` vs `aggregate()`
- ▶ Múltiplas chaves de agrupamento
- ▶ Funções customizadas

Revisão da Aula (cont.)

Bloco 3 - Combinação:

- ▶ `concat()`: empilhamento vertical/horizontal
- ▶ `merge()`: JOIN estilo SQL (inner, left, right, outer)
- ▶ Múltiplas chaves e colunas diferentes
- ▶ `join()`: merge pelo índice
- ▶ Sufixos e indicadores de origem

Bloco 4 - Reshape:

- ▶ `pivot_table()`: criar tabelas de resumo
- ▶ `crosstab()`: tabulação cruzada
- ▶ `melt()`: wide to long format
- ▶ `stack()/unstack()`: reorganizar níveis
- ▶ Tratamento de duplicatas

Conceitos-Chave para Lembrar

1. **Prefira vetorização** a `apply()` quando possível
2. **str accessor** para operações com strings
3. **GroupBy = Split-Apply-Combine**
4. **transform()** mantém shape, `agg()` reduz
5. **merge()** para relacionamentos, `concat()` para empilhamento
6. **Inner join** é padrão no `merge()`
7. **pivot_table()** cria resumos agregados
8. **melt()** converte wide → long
9. **Method chaining** torna pipelines legíveis

Performance: Dicas Importantes

Otimização de código Pandas:

1. Vetorização } `apply()` > loops

- ▶ Use operações vetorizadas sempre que possível
- ▶ `apply()` apenas quando necessário
- ▶ NUNCA use loops explícitos

2. Categorical para colunas repetitivas

- ▶ Economiza memória significativamente
- ▶ Acelera operações de grupo

3. Use dtypes apropriados

- ▶ `int32` vs `int64`, `float32` vs `float64`
- ▶ Pode reduzir uso de memória pela metade

4. Filter antes de `groupby`

- ▶ Reduza dados antes de agrupar
- ▶ Menos dados = operações mais rápidas

Workflow Típico de Manipulação

Pipeline padrão:

1. **Carregar:** `read_csv()`, `read_excel()`
2. **Limpar:** `drop_duplicates()`, `dropna()` ou `fillna()`
3. **Transformar:** criar features, `apply()`, `map()`
4. **Enriquecer:** `merge()` com dados adicionais
5. **Agrupar:** `groupby()` para análises agregadas
6. **Pivotar:** `pivot_table()` para resumos
7. **Visualizar:** (próxima aula)
8. **Exportar:** `to_csv()`, `to_excel()`

💡 Nota Importante

Use method chaining para pipelines mais legíveis

Ao manipular dados:

1. Documente transformações

- ▶ Comente o "porquê", não apenas o "o quê"
- ▶ Use nomes descritivos de variáveis

2. Verifique resultados intermediários

- ▶ Use head(), shape, info() frequentemente
- ▶ Confirme que transformações funcionaram

3. Mantenha dados originais

- ▶ Trabalhe com cópias: `df_work = df.copy()`
- ▶ Não modifique dados brutos

4. Valide agregações

- ▶ Verifique totais e contagens
- ▶ Confirme que groupby não perdeu dados

Erros Comuns e Como Evitar

1. GroupBy sem agregação

- ▶ `df.groupby('col')` apenas cria objeto
- ▶ Precisa de `.agg()`, `.mean()`, etc.

2. Merge sem especificar 'on'

- ▶ Pode usar colunas erradas automaticamente
- ▶ Sempre especifique explicitamente

3. Confundir concat com merge

- ▶ `concat` = empilhar, `merge` = JOIN
- ▶ Use cada um para o propósito certo

4. Esquecer de resetar índice após operações

- ▶ `groupby`, `sort_values` deixam índice bagunçado
- ▶ Use `reset_index()` quando necessário

Próxima Aula: EDA - Parte 1

Aula 07 - O que vem:

- ▶ Estatística descritiva completa
- ▶ Correlação e relações entre variáveis
- ▶ Identificação de outliers
- ▶ Visualização com Matplotlib
- ▶ Múltiplos gráficos e layouts
- ▶ Dashboard simples

Preparação:

- ▶ Revisar conceitos de estatística básica
- ▶ Praticar transformações desta aula
- ▶ Explorar Iris dataset (próxima aula)
- ▶ Pensar em perguntas analíticas

Dicas de Estudo

Para dominar manipulação de dados:

1. **Pratique com datasets diversos**

- ▶ Cada dataset tem desafios únicos
- ▶ Kaggle, UCI ML Repository

2. **Crie seus próprios pipelines**

- ▶ Do início ao fim: carregar → transformar → analisar
- ▶ Documente cada etapa

3. **Estude exemplos reais**

- ▶ Notebooks no Kaggle
- ▶ Projetos open source

4. **Domine groupby**

- ▶ É uma das operações mais poderosas
- ▶ Pratique diferentes agregações

Recursos Adicionais

Documentação:

- ▶ Pandas User Guide: Group By
- ▶ Pandas Cookbook
- ▶ Stack Overflow: tag [pandas]

Tutoriais:

- ▶ Real Python: Pandas GroupBy
- ▶ Kaggle Learn: Data Cleaning
- ▶ DataCamp: Merging DataFrames

Prática:

- ▶ 100 pandas puzzles (GitHub)
- ▶ Exercícios do Kaggle
- ▶ Projetos pessoais

Preparação para Entrega Quinzenal

Segunda entrega (31/10) - Status:

- ▶ **Aula 5:** Pandas básico ✓
- ▶ **Aula 6:** Manipulação ✓ (hoje)
- ▶ **Aula 7:** EDA Parte 1 (próxima)
- ▶ **Aula 8:** EDA Parte 2

Tempo de prática em aula:

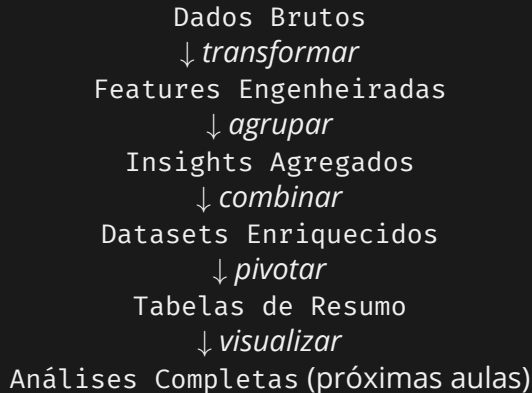
- ▶ Semana 6: 1 hora (após Aula 7)
- ▶ Semana 8: 1 hora (após Aula 8) + entrega

Lembrete:

- ▶ Trabalho envolve Titanic E Wine Quality
- ▶ Feature engineering será importante
- ▶ Comece a pensar nas análises

Resumo: Poder do Pandas

O que conseguimos fazer agora:



Exercício Prático

Tempo: 60 minutos

Entrega: via Moodle (notebook)

Dataset: titanic.csv

Tarefas:

1. (atualizado durante a aula)

Notebook: Disponível no Moodle

Obrigado!

Próxima aula: Análise Exploratória de Dados - Parte 1
Terça-feira, 28/10

Dúvidas: via Moodle ou atendimento