



Pandas in Python

Managing Datasets

What We'll Cover Today

01

Data Organization

Master sorting techniques and eliminate duplicate records from your datasets

03

Data Aggregation

Harness GroupBy operations to summarize and analyze data patterns

02

Data Integration

Combine multiple datasets using various merge strategies and join operations

04

Column Operations

Transform and manipulate columns using pandas' powerful built-in functions

20883	29.13	2009	16	-	-	7	-	-	-
21818	20.41	20010	19	26	023	91	80	53	1
1198	2021	2298	11	50	55	89	-	-	-
2188	29.10	21.300	48	98	-	-	-	-	-
2188	21.02	06.503	38	241	19	8	-	-	-
2556	98.06	20.929	1	28	38	9	-	-	-
11984	21.36	1224	20	0	4	-	-	-	-
2100	20.23	2220	10	3632	59	54	18	62	-
2008	20818	2108	1	38	11	36	6	-	-
2198	21.96	1990	1	21	0	0	-	-	-
8010	8136	23506	102	98	38	29	-	-	-

Data Cleaning Fundamentals

Before analyzing data, we need to ensure it's clean, organized, and ready for insights. Let's start with the essential operations that form the foundation of every data analysis workflow.

Sorting Data Like a Pro

Why Sorting Matters

Proper sorting helps you:

- Identify patterns and outliers quickly
- Prepare data for efficient merging
- Improve readability of your datasets
- Optimize performance for certain operations



Basic Sorting Operations

```
# Sort by single column  
df.sort_values('column_name')
```

```
# Sort by multiple columns  
df.sort_values(['col1', 'col2'])
```

```
# Sort in descending order  
df.sort_values('column_name', ascending=False)
```

```
# Sort by index  
df.sort_index()
```

The `sort_values()` method is your go-to tool for organizing data. Use `ascending=False` for descending order, and pass a list of column names to sort by multiple criteria with priority order.

Advanced Sorting Techniques

Custom Sort Orders

```
# Define custom category order
custom_order = ['Low', 'Medium', 'High']
df['category'] = pd.Categorical(
    df['category'],
    categories=custom_order,
    ordered=True
)
df.sort_values('category')
```

Handling Missing Values

```
# Control where NaN values appear
df.sort_values('column', na_position='first')
df.sort_values('column', na_position='last')
```

In-Place vs. Copy Operations

Understanding the Difference

```
# Creates a new DataFrame (default)  
sorted_df = df.sort_values('column')
```

```
# Modifies original DataFrame  
df.sort_values('column', inplace=True)
```

```
# Always safer to create copies first  
df_copy = df.copy()  
df_copy.sort_values('column', inplace=True)
```

Use `inplace=True` carefully - it permanently modifies your original data. Creating copies gives you flexibility to experiment without losing your original dataset.



Eliminating Duplicates

Duplicate records can skew your analysis and waste computational resources. Let's explore pandas' powerful deduplication capabilities to keep your data clean and accurate.



Finding Duplicates

```
# Check for duplicate rows  
df.duplicated()  
  
# Check for duplicates in specific columns  
df.duplicated(subset=['col1', 'col2'])  
  
# View actual duplicate rows  
df[df.duplicated()]  
  
# Count total duplicates  
df.duplicated().sum()
```

The `duplicated()` method returns a boolean Series indicating duplicate rows. By default, it marks all duplicates except the first occurrence as True.

Removing Duplicates Strategically

1

Keep First Occurrence

```
df.drop_duplicates(keep='first')
```

Default behavior - keeps the first instance of each duplicate group

2

Keep Last Occurrence

```
df.drop_duplicates(keep='last')
```

Useful when later entries contain more recent or complete information

3

Remove All Duplicates

```
df.drop_duplicates(keep=False)
```

Removes all instances of duplicated rows, keeping only unique entries

Selective Deduplication

Target Specific Columns

```
# Remove duplicates based on specific columns  
df.drop_duplicates(subset=['name', 'email'])
```

```
# Multiple column deduplication
```

```
df.drop_duplicates(  
    subset=['customer_id', 'product_id'],  
    keep='last'  
)
```

Real-World Example

Customer databases often have multiple entries for the same person. Use `subset` to identify duplicates based on key identifying fields like email or customer ID, while preserving other valuable data differences.

Duplicate Detection Best Practices



Investigate First

Always examine your duplicates before removing them. Use `df[df.duplicated()].head()` to understand what you're dealing with.



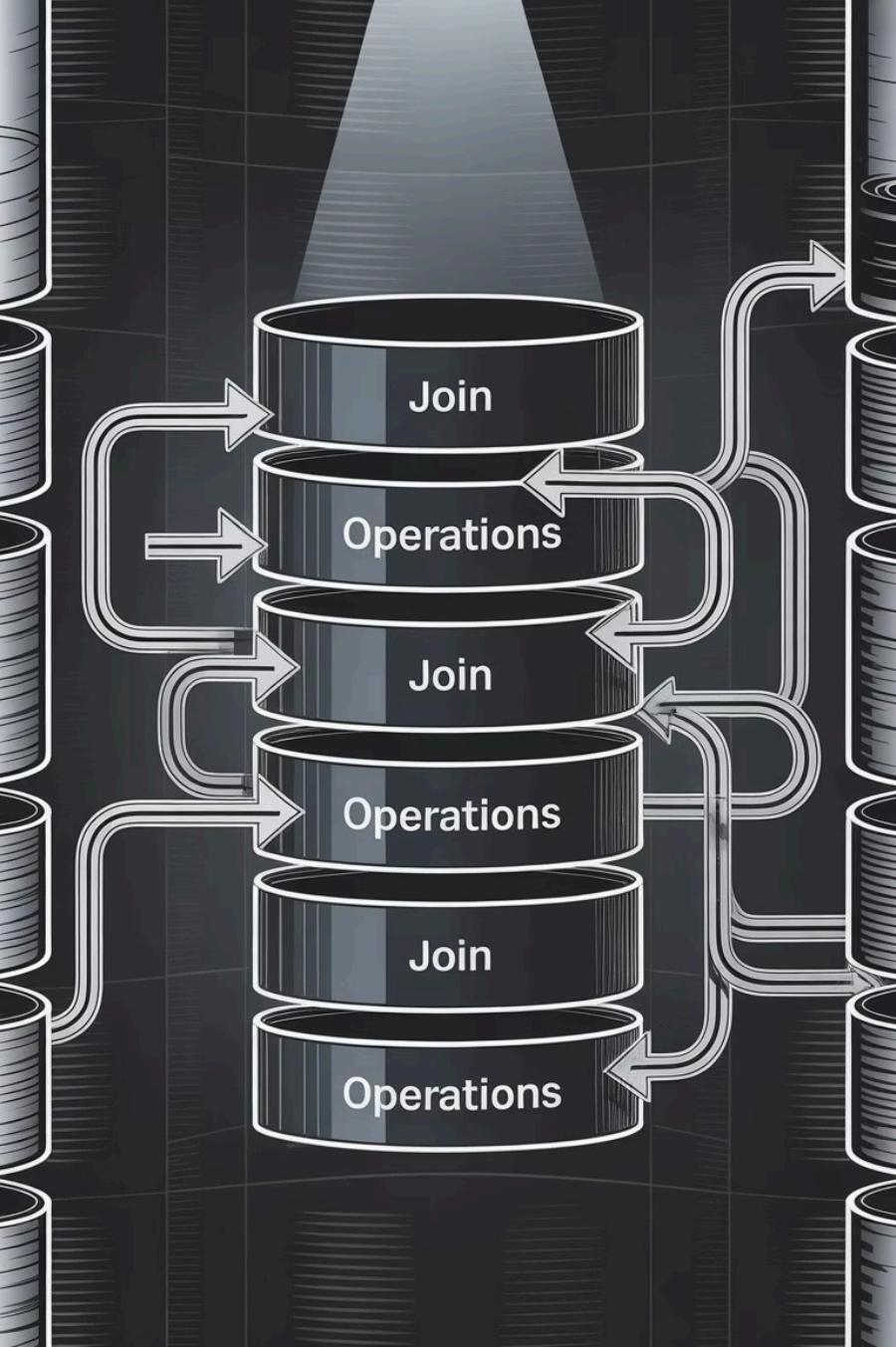
Backup Your Data

Create copies before deduplication. Once you drop duplicates, the original structure is lost unless you've saved it.



Be Strategic

Consider business logic when choosing `keep='first'`, `'last'`, or `False`. The context matters more than the technical implementation.

A vertical stack of five cylindrical containers representing datasets. From top to bottom, they are labeled: 'Join', 'Operations', 'Join', 'Operations', and 'Join'. Arrows show data flowing from the bottom 'Operations' container up through the stack, with 'Join' operations occurring between each consecutive pair of containers. The background features a dark grid pattern.

Merging Datasets

Combining data from multiple sources is essential for comprehensive analysis. Master pandas' merge capabilities to unlock insights from related datasets.

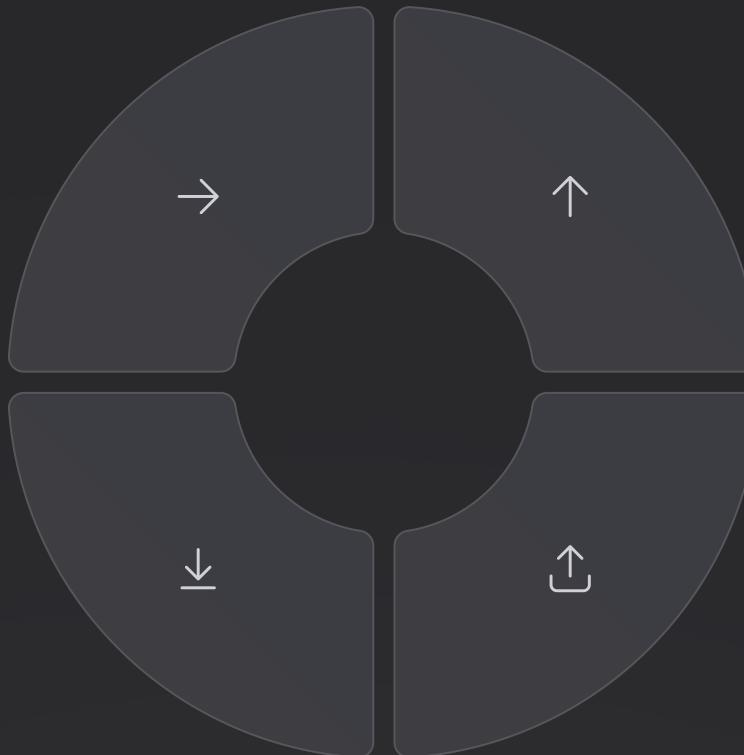
Understanding Join Types

Inner Join

Returns only matching records from both datasets. Most restrictive but ensures data consistency.

Outer Join

Includes all records from both datasets. Creates the most comprehensive view with NaN for missing matches.



Left Join

Keeps all records from left dataset, adds matching records from right. Preserves your primary dataset.

Right Join

Keeps all records from right dataset, adds matching records from left. Less commonly used than left join.

Basic Merge Syntax

```
# Basic merge on common column  
pd.merge(df1, df2, on='common_column')  
  
# Specify join type  
pd.merge(df1, df2, on='id', how='left')  
  
# Merge on multiple columns  
pd.merge(df1, df2, on=['col1', 'col2'])  
  
# Different column names  
pd.merge(df1, df2, left_on='id1', right_on='id2')
```

The `pd.merge()` function is pandas' primary tool for combining DataFrames. Always specify the join type with `how` parameter for clarity and predictable results.

Handling Merge Conflicts

Column Name Conflicts

```
# Automatic suffixes  
pd.merge(df1, df2, on='id', suffixes=('_left', '_right'))
```

```
# Custom suffixes  
pd.merge(df1, df2, on='id', suffixes=('_main', '_lookup'))
```

Missing Key Columns

```
# Explicit column mapping  
pd.merge(df1, df2,  
        left_on='customer_id',  
        right_on='cust_id',  
        how='inner')
```

Index-Based Merging

```
# Merge on index  
pd.merge(df1, df2, left_index=True,  
        right_index=True)
```

```
# Mix index and columns  
pd.merge(df1, df2, left_on='id',  
        right_index=True)
```

Advanced Merge Techniques

Validation and Indicators

```
# Validate merge expectations  
pd.merge(df1, df2,  
        on='id',  
        validate='one_to_many')
```

```
# Add merge indicator  
pd.merge(df1, df2,  
        on='id',  
        how='outer',  
        indicator=True)
```

Why This Matters

Use `validate` to catch data quality issues early. The `indicator` parameter adds a special column showing the source of each row: 'left_only', 'right_only', or 'both'.



Concatenation vs. Merging

Use pd.concat() When

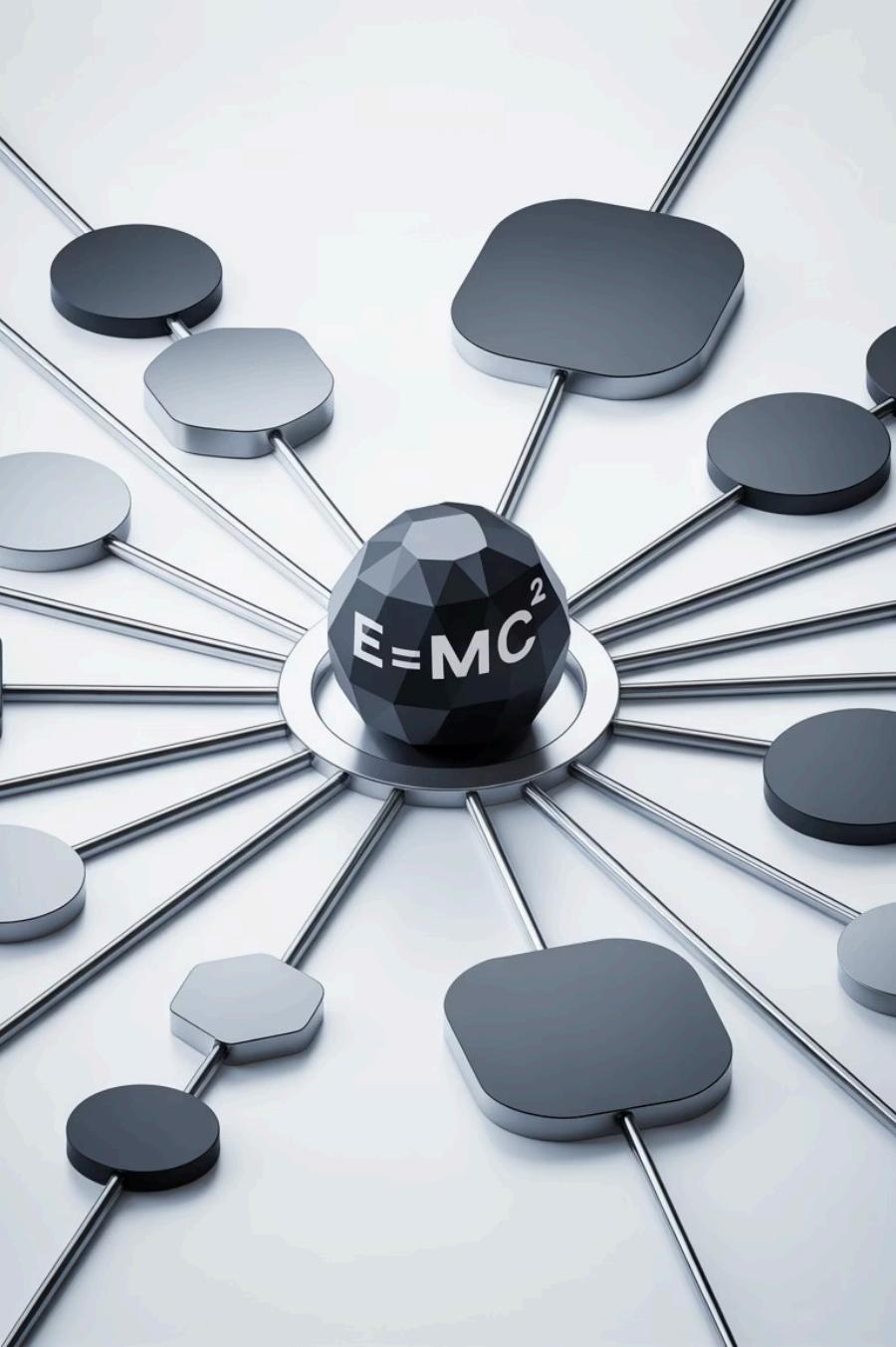
- Stacking DataFrames vertically
- Same column structure
- Appending new rows
- Simple data combination

```
pd.concat([df1, df2], ignore_index=True)
```

Use pd.merge() When

- Joining based on key columns
- Different column structures
- Relational database-style joins
- Complex matching logic needed

```
pd.merge(df1, df2, on='key_column')
```



GroupBy Operations

GroupBy is pandas' most powerful feature for data aggregation and analysis. Think of it as SQL's GROUP BY clause, but with much more flexibility and functionality.

The GroupBy Philosophy

Split

Divide the DataFrame into groups based on specified criteria or column values

Apply

Perform operations (aggregation, transformation, filtration) on each group independently

Combine

Merge the results back into a single DataFrame with meaningful structure

Basic GroupBy Operations

```
# Group by single column  
df.groupby('category').sum()  
  
# Group by multiple columns  
df.groupby(['region', 'category']).mean()  
  
# Access specific columns after grouping  
df.groupby('category')['sales'].sum()  
  
# Multiple aggregations  
df.groupby('category').agg(['sum', 'mean', 'count'])
```

Start simple with basic aggregations like `sum()`, `mean()`, and `count()`. These form the foundation for more complex analysis patterns.

Essential Aggregation Functions



sum()

Total values in each group. Perfect for sales totals, quantities, or any cumulative metrics.



mean()

Average values per group. Essential for understanding typical performance or behavior patterns.



#

count()

Number of non-null entries. Reveals data completeness and group sizes for analysis validity.

max() / min()

Extreme values in each group. Identify outliers, ranges, and performance boundaries.



()

std()

Standard deviation shows variability. High values indicate inconsistent performance within groups.

size()

Count of all entries including nulls. Different from count() - shows true group membership.

Advanced Aggregation with agg()

Multiple Functions per Column

```
# Different functions for each column
df.groupby('category').agg({
    'sales': ['sum', 'mean'],
    'quantity': 'count',
    'price': ['min', 'max']
})
```

Custom Aggregation Functions

```
# Define custom function
def price_range(series):
    return series.max() - series.min()

# Apply custom function
df.groupby('category').agg({
    'price': price_range,
    'sales': 'sum'
})
```

Transform and Filter Operations

1

Transform

Apply functions that return same-sized results. Perfect for standardizing data within groups or calculating group-relative metrics.

```
# Standardize within groups
df['sales_zscore'] = df.groupby('region')['sales'].transform(
    lambda x: (x - x.mean()) / x.std()
)
```

2

Filter

Keep only groups that meet specific criteria. Useful for focusing analysis on significant or qualifying segments.

```
# Keep groups with more than 10 entries
large_groups = df.groupby('category').filter(
    lambda x: len(x) > 10
)
```

Working with Multi-Level Indexes

```
# Create multi-level grouping  
grouped = df.groupby(['region', 'category']).sum()  
  
# Reset index to flatten  
grouped.reset_index()  
  
# Access specific levels  
grouped.loc['North']  
grouped.loc[('North', 'Electronics')]  
  
# Unstack to create pivot-like structure  
grouped.unstack('category')
```

Multi-level grouping creates hierarchical indexes. Use `reset_index()` to convert back to regular columns, or `unstack()` to create pivot-table-like views for easier analysis.

Column Manipulation

Effective column manipulation is the backbone of data preprocessing. Master these techniques to reshape, transform, and enhance your datasets for analysis.



Creating and Modifying Columns

Basic Column Creation

```
# Simple assignment  
df['new_column'] = value  
  
# From existing columns  
df['total'] = df['price'] * df['quantity']  
  
# Conditional creation  
df['category'] = df['score'].apply(  
    lambda x: 'High' if x > 80 else 'Low'  
)
```

Multiple Column Operations

```
# Assign multiple columns at once  
df = df.assign(  
    profit=df['revenue'] - df['cost'],  
    margin=df['profit'] / df['revenue']  
)  
  
# Using eval for complex expressions  
df.eval('profit_margin = (revenue - cost) / revenue',  
    inplace=True)
```

Essential Column Functions

String Operations

```
# String methods  
df['name'].str.lower()  
df['name'].str.contains('pattern')  
df['name'].str.replace('old', 'new')  
df['name'].str.split(',')  
  
# Extract patterns  
df['phone'].str.extract(r'(\d{3})-(\d{3})-(\d{4})')
```

Numeric Operations

```
# Mathematical functions  
df['value'].abs()  
df['value'].round(2)  
df['value'].clip(lower=0, upper=100)  
  
# Statistical functions  
df['value'].rank()  
df['value'].pct_change()  
df['value'].cumsum()
```

Advanced Column Transformations

apply() Function

Most versatile transformation tool. Apply any function to each element in a column or across rows/columns.

```
df['column'].apply(custom_function)
```

where() Function

Conditional replacement. Keep values that meet condition, replace others with specified value or NaN.

```
df['score'].where(df['score'] > 50, 0)
```

1

2

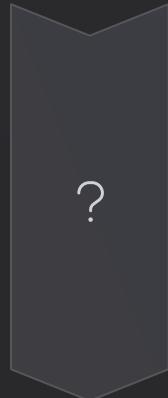
3

map() Function

Perfect for value replacement using dictionaries. More efficient than apply() for simple mappings.

```
df['grade'].map({'A': 90, 'B': 80, 'C': 70})
```

Handling Missing Data in Columns



Detect Missing Values

```
df.isnull().sum() # Count nulls per column  
df.info()       # Overview of data types and nulls  
df.describe()   # Statistical summary
```



Fill Missing Values

```
df.fillna(value)        # Fill with specific value  
df.fillna(method='ffill') # Forward fill  
df.fillna(df.mean())    # Fill with column mean
```



Remove Missing Data

```
df.dropna() # Drop rows with any nulls  
df.dropna(subset=['column']) # Drop based on specific columns  
df.dropna(thresh=2) # Keep rows with at least 2 non-null values
```

Data Type Conversions

Common Conversions

```
# Numeric conversions  
df['column'].astype(int)  
df['column'].astype(float)  
  
# String conversions  
df['column'].astype(str)  
  
# DateTime conversions  
pd.to_datetime(df['date_column'])  
  
# Categorical data  
df['category'].astype('category')
```

Error Handling

```
# Handle conversion errors  
pd.to_numeric(df['column'], errors='coerce')  
  
# Check data types  
df.dtypes  
df.info()  
  
# Memory optimization  
df['category'] = df['category'].astype('category')
```

Always verify data types after import. Proper types improve performance and prevent analysis errors.

Performance Tips for Large Datasets



Use Vectorized Operations

Avoid loops! Use pandas vectorized functions which are implemented in C and much faster than Python loops.

```
df['result'] = df['col1'] + df['col2'] #  
Fast  
# Not: df['result'] = df.apply(lambda x:  
x['col1'] + x['col2'], axis=1)
```



Optimize Data Types

Use appropriate data types to reduce memory usage. Categorical data and smaller numeric types can significantly improve performance.

```
df['category'] =  
df['category'].astype('category')  
df['small_int'] =  
df['big_int'].astype('int8')
```



Filter Early

Apply filters as early as possible in your data pipeline to reduce the amount of data being processed.

```
# Filter first, then process  
df_filtered = df[df['score'] > 50]  
results =  
df_filtered.groupby('category').mean()
```

Putting It All Together: Real-World Example

```
# Complete data processing pipeline
def process_sales_data(df):
    # 1. Clean and deduplicate
    df_clean = df.drop_duplicates(subset=['order_id'])

    # 2. Handle missing values
    df_clean['customer_segment'] = df_clean['customer_segment'].fillna('Unknown')

    # 3. Create derived columns
    df_clean['profit'] = df_clean['revenue'] - df_clean['cost']
    df_clean['profit_margin'] = df_clean['profit'] / df_clean['revenue']

    # 4. Group and aggregate
    summary = df_clean.groupby(['region', 'product_category']).agg({
        'revenue': 'sum',
        'profit': 'sum',
        'order_id': 'count'
    }).rename(columns={'order_id': 'order_count'})

    # 5. Sort and return
    return summary.sort_values('revenue', ascending=False)
```

Best Practices Summary

1 Always Explore First

Use `df.info()`, `df.describe()`, and `df.head()` to understand your data before applying operations.

3 Validate Your Results

Always check your output. Use `.shape`, `.dtypes`, and sample the results to ensure operations worked as expected.

2 Chain Operations Wisely

Method chaining improves readability, but break long chains for debugging. Use parentheses to span multiple lines.

4 Document Your Process

Comment complex operations and document your assumptions. Future you (and your colleagues) will appreciate the clarity.



Next Steps in Your Pandas Journey



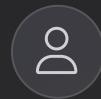
Practice with Real Data

Apply these techniques to your own datasets. Start with small, manageable files and gradually work with larger, more complex data.



Explore Advanced Features

Learn about pivot tables, time series analysis, and advanced indexing. These build on the fundamentals you've mastered today.



Join the Community

Engage with pandas communities, read documentation regularly, and share your experiences. The pandas ecosystem is constantly evolving.

You now have the essential tools to clean, merge, group, and manipulate data effectively. These operations form the foundation of professional data analysis workflows.