

# Pandas Mastery

From Beginner to Advanced

A Comprehensive Guide to Data Analysis with Python

Anshuman Singh

November 27, 2025

© 2025 Anshuman Singh

All rights reserved.

This work may be distributed and/or modified under the conditions  
of the Creative Commons Attribution 4.0 International License.

First Edition: November 27, 2025

ISBN: 978-0-000000-00-0

Cover Design: Data Science Edition

Typeset in L<sup>A</sup>T<sub>E</sub>X

# Contents

<b>I</b>	<b>Getting Started with Pandas</b>	<b>7</b>
<b>1</b>	<b>Introduction to Pandas</b>	<b>8</b>
1.1	What is Pandas?	8
1.1.1	Why Use Pandas?	8
1.2	Installing Pandas	9
1.2.1	Verifying Installation	9
1.3	Importing Pandas	9
1.3.1	Why the pd alias?	9
1.4	Pandas Ecosystem	10
1.5	Exercises	10
1.6	Solutions	10
<b>2</b>	<b>Series and DataFrames: Core Data Structures</b>	<b>12</b>
2.1	Understanding Series	12
2.1.1	Creating Series Objects	12
2.1.2	Series Attributes and Methods	12
2.2	Understanding DataFrames	13
2.2.1	Creating DataFrames	13
2.3	Basic DataFrame Operations	14
2.4	DataFrame Attributes and Metadata	14
2.5	Exercises	15
2.6	Solutions	15
<b>3</b>	<b>Indexing, Slicing, and Data Selection</b>	<b>17</b>
3.1	Understanding Indexing Methods	17
3.1.1	loc[] - Label-based Indexing	17
3.1.2	iloc[] - Integer-based Indexing	17
3.1.3	Boolean Indexing	18
3.2	Setting Values with Indexing	18
3.3	Exercises	19
3.4	Solutions	19
<b>4</b>	<b>Data Cleaning and Preparation</b>	<b>20</b>
4.1	Handling Missing Data	20
4.1.1	Detecting Missing Values	20
4.1.2	Removing Missing Data	20
4.1.3	Filling Missing Data	21
4.2	Data Type Conversion	21
4.3	String Operations	22
4.4	Exercises	22
4.5	Solutions	23

<b>II Intermediate Pandas: Data Manipulation</b>	<b>24</b>
<b>5 Data Input/Output Operations</b>	<b>25</b>
5.1 Reading Different File Formats	25
5.2 Writing Data to Files	25
5.3 File Reading Parameters and Options	26
5.4 Handling Large Files	26
5.5 Exercises	27
5.6 Solutions	27
<b>6 Data Aggregation and Group Operations</b>	<b>30</b>
6.1 GroupBy Fundamentals	30
6.2 Aggregation Methods	30
6.3 Transform and Filter Operations	31
6.4 Pivot Tables	31
6.5 Cross Tabulation	32
6.6 Exercises	32
6.7 Solutions	33
<b>7 Data Merging, Joining, and Concatenating</b>	<b>35</b>
7.1 Concatenation Basics	35
7.2 Merging DataFrames	35
7.3 Advanced Merging Techniques	36
7.4 Joining DataFrames	37
7.5 Exercises	38
7.6 Solutions	38
<b>III Advanced Pandas: Optimization and Real-World Applications</b>	<b>41</b>
<b>8 Advanced Data Manipulation with MultiIndex</b>	<b>42</b>
8.1 Creating and Understanding MultiIndex	42
8.2 Indexing with MultiIndex	42
8.3 Reshaping with Stack and Unstack	43
8.4 Advanced MultiIndex Operations	43
8.5 Exercises	44
8.6 Solutions	44
<b>9 Performance Optimization and Vectorization</b>	<b>46</b>
9.1 Understanding Vectorization	46
9.2 Memory Optimization	46
9.3 Efficient Data Processing	47
9.4 Working with Large Datasets	48
9.5 Exercises	49
9.6 Solutions	49
<b>10 Time Series Analysis with Pandas</b>	<b>53</b>
10.1 Working with DateTime Data	53
10.2 Time-Based Indexing and Resampling	53
10.3 Rolling Windows and Expanding Operations	54

10.4	Time Series Decomposition . . . . .	54
10.5	Lag Features and Time Shifts . . . . .	55
10.6	Exercises . . . . .	55
10.7	Solutions . . . . .	56
<b>11</b>	<b>Working with Text Data . . . . .</b>	<b>58</b>
11.1	Basic String Operations . . . . .	58
11.2	String Methods and Pattern Matching . . . . .	58
11.3	Text Analysis and NLP Basics . . . . .	59
11.4	Regular Expressions with Pandas . . . . .	60
11.5	Categorical Text Analysis . . . . .	60
11.6	Exercises . . . . .	61
11.7	Solutions . . . . .	61
<b>12</b>	<b>Advanced Visualization with Pandas . . . . .</b>	<b>65</b>
12.1	Basic Plotting with Pandas . . . . .	65
12.2	Advanced Plot Types . . . . .	66
12.3	Time Series Visualization . . . . .	67
12.4	Interactive Visualizations . . . . .	67
12.5	Custom Styling and Themes . . . . .	69
12.6	Exercises . . . . .	69
12.7	Solutions . . . . .	70
<b>13</b>	<b>Real-World Case Studies and Best Practices . . . . .</b>	<b>74</b>
13.1	Complete Data Analysis Workflow . . . . .	74
13.2	Performance Optimization Case Study . . . . .	77
13.3	Best Practices and Code Quality . . . . .	80
13.4	Final Project: Complete Business Analysis . . . . .	84
13.5	Exercises . . . . .	89
13.6	Solutions . . . . .	89
<b>14</b>	<b>Working with Different Data Sources . . . . .</b>	<b>94</b>
14.1	Database Operations with Pandas . . . . .	94
14.2	Web APIs and JSON Data . . . . .	95
14.3	Web Scraping with Pandas . . . . .	97
14.4	Working with Excel and Office Documents . . . . .	98
14.5	Cloud Storage and Big Data . . . . .	99
14.6	Real-time Data and Streaming . . . . .	101
14.7	Exercises . . . . .	102
14.8	Solutions . . . . .	103
	<b>Glossary . . . . .</b>	<b>109</b>
	<b>About the Author . . . . .</b>	<b>111</b>

# List of Figures

- 1.1 Pandas data structures: Series (1-dimensional) and DataFrame (2-dimensional) . 8

Anshuman Singh — [anshuman365.github.io](https://github.com/anshuman365)

# Part I

## Getting Started with Pandas

Anshuman Singh — anshuman365@gmail.com

## Introduction to Pandas

### 1.1 What is Pandas?

Pandas is a fast, powerful, flexible, and easy-to-use open-source data analysis and manipulation tool built on top of the Python programming language. The name "Pandas" is derived from the term "panel data", an econometrics term for multidimensional structured data sets.

#### 1.1.1 Why Use Pandas?

Pandas has become the de facto standard for data manipulation in Python due to its intuitive data structures and powerful data analysis capabilities. It bridges the gap between spreadsheet applications and programming languages, providing both ease of use and programmatic power.

Key features of Pandas include:

- **DataFrame object:** For data manipulation with integrated indexing
- **Data I/O:** Tools for reading and writing data between in-memory data structures and different file formats (CSV, Excel, JSON, SQL, etc.)
- **Data alignment:** Integrated handling of missing data with intelligent data alignment
- **Reshaping and pivoting:** Powerful tools for reshaping and pivoting data sets
- **Label-based slicing:** Sophisticated indexing and subsetting of large data sets
- **Group by operations:** Split-apply-combine operations on data sets
- **High performance merging:** Efficient merging and joining of data
- **Time series functionality:** Date range generation and frequency conversion, moving window statistics, date shifting and lagging

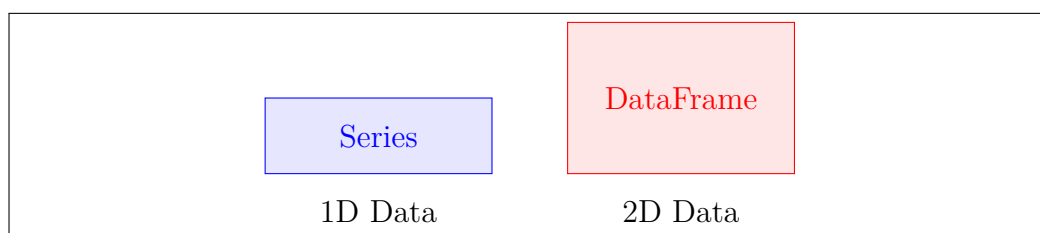


Figure 1.1: Pandas data structures: Series (1-dimensional) and DataFrame (2-dimensional)



## 1.2 Installing Pandas

Pandas can be installed using various package managers. Here are the most common methods:

```
1 # Using pip (most common)
2 pip install pandas
3
4 # Using conda (for Anaconda/Miniconda users)
5 conda install pandas
6
7 # For specific version
8 pip install pandas==2.0.0
9
10 # Installing with optional dependencies
11 pip install pandas[performance]
12
13 # From source (advanced users)
14 git clone https://github.com/pandas-dev/pandas.git
15 cd pandas
16 python setup.py install
```

Listing 1.1: Installing Pandas using different methods

### 1.2.1 Verifying Installation

After installation, verify that pandas is properly installed:

```
1 import pandas as pd
2 print(f"Pandas version: {pd.__version__}")
3 print(f"Pandas location: {pd.__file__}")
```

Listing 1.2: Verifying pandas installation

## 1.3 Importing Pandas

The conventional way to import pandas in Python follows these best practices:

```
1 import pandas as pd
2 import numpy as np # Often used together with pandas
3
4 # Common aliases for data visualization
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
8 # Display settings for better output
9 pd.set_option('display.max_columns', None)
10 pd.set_option('display.width', 1000)
```

Listing 1.3: Importing pandas and related libraries

### 1.3.1 Why the pd alias?

The pd alias has become a community standard for several reasons:

- Shorter and more convenient to type
- Consistent with community conventions

- Avoids namespace conflicts
- Makes code more readable and professional

## 1.4 Pandas Ecosystem

Pandas integrates seamlessly with other Python libraries:

- **NumPy**: Foundation for pandas data structures
- **Matplotlib/Seaborn**: For data visualization
- **Scikit-learn**: For machine learning
- **Statsmodels**: For statistical analysis
- **SQLAlchemy**: For database operations

## 1.5 Exercises

1. Install pandas on your system using your preferred method and verify the installation.
2. Create a Python script that imports pandas and prints the version information.
3. What are the main advantages of using pandas over plain Python lists or NumPy arrays for data analysis?
4. List three common file formats that pandas can read and write.

## 1.6 Solutions

```
1.
1 # Verify installation and version
2 import pandas as pd
3 print(f"Pandas version: {pd.__version__}")
4
5 # Check basic functionality
6 data = {'Name': ['Alice', 'Bob'], 'Age': [25, 30]}
7 df = pd.DataFrame(data)
8 print(df)
```

```
1 import pandas as pd
2 import numpy as np
3
4 print(f"Pandas version: {pd.__version__}")
5 print(f"NumPy version: {np.__version__}")
6
7 # Additional system information
8 import sys
9 print(f"Python version: {sys.version}")
```

3. Main advantages of pandas:

- Label-based indexing for intuitive data access

- Built-in handling of missing data
- Powerful group by operations
- Integrated time series functionality
- Easy data I/O with multiple file formats
- Memory efficiency with optimized data structures

4. Three common file formats pandas supports:

- CSV (Comma Separated Values)
- Excel (.xlsx, .xls)
- JSON (JavaScript Object Notation)
- SQL databases
- Parquet
- HDF5

# Chapter 2

## Series and DataFrames: Core Data Structures

### 2.1 Understanding Series

A Series is a one-dimensional labeled array capable of holding any data type (integers, strings, floating point numbers, Python objects, etc.). Think of it as a column in a spreadsheet or a single variable in statistics.

#### 2.1.1 Creating Series Objects

```
1 import pandas as pd
2 import numpy as np
3
4 # Creating Series from list (most common)
5 s1 = pd.Series([1, 3, 5, np.nan, 6, 8])
6 print("Series from list:")
7 print(s1)
8
9 # Creating Series with custom index
10 s2 = pd.Series([10, 20, 30], index=['a', 'b', 'c'])
11 print("\nSeries with custom index:")
12 print(s2)
13
14 # Creating Series from dictionary
15 s3 = pd.Series({'a': 1, 'b': 2, 'c': 3})
16 print("\nSeries from dictionary:")
17 print(s3)
18
19 # Creating Series with specified data type
20 s4 = pd.Series([1, 2, 3], dtype=np.float64)
21 print("\nSeries with specified dtype:")
22 print(s4)
23
24 # Creating Series from NumPy array
25 arr = np.array([1, 2, 3, 4])
26 s5 = pd.Series(arr)
27 print("\nSeries from NumPy array:")
28 print(s5)
```

Listing 2.1: Various ways to create Series

#### 2.1.2 Series Attributes and Methods

```
1 # Create a sample series
2 s = pd.Series([10, 20, 30, 40], index=['a', 'b', 'c', 'd'])
3
4 # Basic attributes
```

```

5 print(f"Values: {s.values}")
6 print(f"Index: {s.index}")
7 print(f"Data type: {s.dtype}")
8 print(f"Shape: {s.shape}")
9 print(f"Size: {s.size}")
10 print(f"Name: {s.name}")
11
12 # Basic methods
13 print(f"Sum: {s.sum()}")
14 print(f"Mean: {s.mean()}")
15 print(f"Standard deviation: {s.std()}")
16 print(f"Minimum: {s.min()}")
17 print(f"Maximum: {s.max()}")
18
19 # Boolean operations
20 print(f"Greater than 25: {s > 25}")

```

Listing 2.2: Working with Series attributes

## 2.2 Understanding DataFrames

A DataFrame is a two-dimensional labeled data structure with columns of potentially different types, similar to a spreadsheet or SQL table. It's the most commonly used pandas object.

### 2.2.1 Creating DataFrames

```

1 # Creating DataFrame from dictionary (most common)
2 data = {
3     'Name': ['Alice', 'Bob', 'Charlie', 'David'],
4     'Age': [25, 30, 35, 40],
5     'City': ['New York', 'London', 'Tokyo', 'Paris'],
6     'Salary': [50000, 60000, 70000, 80000]
7 }
8 df = pd.DataFrame(data)
9 print("DataFrame from dictionary:")
10 print(df)
11
12 # Creating DataFrame with custom index
13 df_indexed = pd.DataFrame(data, index=['EMP001', 'EMP002', 'EMP003', 'EMP004'])
14 print("\nDataFrame with custom index:")
15 print(df_indexed)
16
17 # Creating DataFrame from list of lists
18 data_list = [
19     ['Alice', 25, 'New York', 50000],
20     ['Bob', 30, 'London', 60000],
21     ['Charlie', 35, 'Tokyo', 70000],
22     ['David', 40, 'Paris', 80000]
23 ]
24 columns = ['Name', 'Age', 'City', 'Salary']
25 df_from_list = pd.DataFrame(data_list, columns=columns)
26 print("\nDataFrame from list of lists:")
27 print(df_from_list)
28
29 # Creating DataFrame from NumPy array
30 np_array = np.random.rand(4, 3)
31 df_from_np = pd.DataFrame(np_array, columns=['A', 'B', 'C'])

```

```
32 print("\nDataFrame from NumPy array:")
33 print(df_from_np)
```

Listing 2.3: Various ways to create DataFrames

## 2.3 Basic DataFrame Operations

```
1 # Create sample DataFrame
2 data = {
3     'Product': ['Laptop', 'Mouse', 'Keyboard', 'Monitor', 'Headphones'],
4     'Price': [1000, 25, 75, 300, 150],
5     'Quantity': [5, 20, 15, 8, 12],
6     'Category': ['Electronics', 'Accessories', 'Accessories', 'Electronics', '
7                 Audio']
8 }
9 df = pd.DataFrame(data)
10
11 # Display basic information
12 print("DataFrame head:")
13 print(df.head())
14
15 print("\nDataFrame tail:")
16 print(df.tail(2))
17
18 print("\nDataFrame info:")
19 print(df.info())
20
21 print("\nDataFrame description:")
22 print(df.describe())
23
24 # Shape and dimensions
25 print(f"\nShape: {df.shape}")
26 print(f"Columns: {df.columns.tolist()}")
27 print(f"Index: {df.index.tolist()}")
28
29 # Access columns
30 print("\nSingle column:")
31 print(df['Product'])
32
33 print("\nMultiple columns:")
34 print(df[['Product', 'Price']])
35
36 # Using dot notation (when column names are valid Python identifiers)
37 print("\nUsing dot notation:")
38 print(df.Price)
```

Listing 2.4: Essential DataFrame operations

## 2.4 DataFrame Attributes and Metadata

```
1 # DataFrame attributes
2 print(f"Columns: {df.columns}")
3 print(f"Index: {df.index}")
4 print(f"Shape: {df.shape}")
5 print(f>Data types: {df.dtypes}")
6 print(f"Memory usage: {df.memory_usage(deep=True).sum()} bytes")
```

```

7
8 # Transposing DataFrame
9 print("\nTransposed DataFrame:")
10 print(df.T)
11
12 # Converting to other formats
13 print("\nAs NumPy array:")
14 print(df.values)
15
16 print("\nAs dictionary:")
17 print(df.to_dict())

```

Listing 2.5: Working with DataFrame metadata

## 2.5 Exercises

1. Create a Series that represents daily temperatures for a week (Monday to Sunday) with appropriate values.
2. Create a DataFrame containing information about 5 books with columns: Title, Author, Year, Pages, and Price.
3. From your books DataFrame, extract only the Title and Price columns.
4. Calculate the average price of books in your DataFrame.
5. Add a new column 'Rating' to your books DataFrame with random ratings between 1 and 5.

## 2.6 Solutions

```

1.
1 temperatures = pd.Series(
2     [22, 24, 23, 25, 26, 27, 26],
3     index=['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', '
4           Saturday', 'Sunday']
5 )
6 print(temperatures)

```

```

1 books_data = {
2     'Title': ['Python Basics', 'Data Science Guide', 'ML Advanced', 'Deep
3               Learning', 'Pandas Cookbook'],
4     'Author': ['Author A', 'Author B', 'Author C', 'Author D', 'Author E'],
5     'Year': [2020, 2021, 2022, 2023, 2024],
6     'Pages': [300, 450, 520, 600, 350],
7     'Price': [29.99, 39.99, 49.99, 59.99, 34.99]
8 }
9 books_df = pd.DataFrame(books_data)
10 print(books_df)

```

```

3.
1 title_price = books_df[['Title', 'Price']]
2 print(title_price)

```

```
1 average_price = books_df['Price'].mean()
2 print(f"Average price: ${average_price:.2f}")
```

```
4.
1 import random
2 books_df['Rating'] = [random.randint(1, 5) for _ in range(len(books_df))]
3 print(books_df)
```



## Indexing, Slicing, and Data Selection

### 3.1 Understanding Indexing Methods

Pandas provides several methods for selecting and slicing data. Understanding these methods is crucial for efficient data manipulation.

#### 3.1.1 `loc[]` - Label-based Indexing

```
1 # Create sample DataFrame
2 data = {
3     'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eva'],
4     'Age': [25, 30, 35, 40, 28],
5     'City': ['NY', 'London', 'Tokyo', 'Paris', 'Berlin'],
6     'Salary': [70000, 80000, 90000, 100000, 75000]
7 }
8 df = pd.DataFrame(data)
9 df.index = ['A', 'B', 'C', 'D', 'E'] # Custom index
10
11 # Single row by label
12 print("Row A:")
13 print(df.loc['A'])
14
15 # Multiple rows by labels
16 print("\nRows A and C:")
17 print(df.loc[['A', 'C']])
18
19 # Rows and specific columns
20 print("\nRows A-C, Name and Salary columns:")
21 print(df.loc['A':'C', ['Name', 'Salary']])
22
23 # Boolean indexing with loc
24 print("\nPeople with Salary > 80000:")
25 print(df.loc[df['Salary'] > 80000])
26
27 # Conditional selection on specific columns
28 print("\nNames of people in London:")
29 print(df.loc[df['City'] == 'London', 'Name'])
```

Listing 3.1: Using `loc` for label-based indexing

#### 3.1.2 `iloc[]` - Integer-based Indexing

```
1 # Integer-based indexing
2 print("First row:")
3 print(df.iloc[0])
4
5 print("\nFirst three rows:")
6 print(df.iloc[0:3])
7
```

```

8 print("\nSpecific rows and columns by position:")
9 print(df.iloc[[0, 2, 4], [0, 2]]) # Rows 0,2,4 and columns 0,2
10
11 print("\nLast two rows:")
12 print(df.iloc[-2:])
13
14 # Using with conditions (less common but possible)
15 print("\nUsing iloc with conditions:")
16 print(df.iloc[(df['Age'] > 30).values])

```

Listing 3.2: Using iloc for integer-based indexing

### 3.1.3 Boolean Indexing

```

1 # Single condition
2 young_employees = df[df['Age'] < 30]
3 print("Employees younger than 30:")
4 print(young_employees)
5
6 # Multiple conditions using & (and), | (or)
7 high_salary_young = df[(df['Salary'] > 80000) & (df['Age'] < 40)]
8 print("\nYoung employees with high salary:")
9 print(high_salary_young)
10
11 # Using isin() for multiple values
12 cities = ['London', 'Paris']
13 in_europe = df[df['City'].isin(cities)]
14 print("\nEmployees in European cities:")
15 print(in_europe)
16
17 # Using str methods for string filtering
18 names_with_a = df[df['Name'].str.contains('a', case=False)]
19 print("\nNames containing 'a':")
20 print(names_with_a)
21
22 # Using query() method (alternative syntax)
23 high_paid = df.query('Salary > 85000 and Age > 30')
24 print("\nHigh paid experienced employees:")
25 print(high_paid)

```

Listing 3.3: Boolean indexing techniques

## 3.2 Setting Values with Indexing

```

1 # Setting single value
2 df.loc['A', 'Salary'] = 75000
3
4 # Setting multiple values
5 df.loc['B':'D', 'Salary'] = [85000, 95000, 105000]
6
7 # Conditional setting
8 df.loc[df['City'] == 'London', 'Salary'] = 90000
9
10 # Using mask to set values
11 df['Bonus'] = 0
12 df.loc[df['Salary'] > 90000, 'Bonus'] = 10000

```

```

13
14 print("Updated DataFrame:")
15 print(df)

```

Listing 3.4: Modifying data using indexing

### 3.3 Exercises

1. Create a DataFrame with 10 rows of sample sales data including: Product, Category, Price, Quantity, Date.
2. Use loc to select all products in a specific category with price greater than 50.
3. Use iloc to select the first 5 rows and only the Product and Price columns.
4. Use boolean indexing to find products with quantity less than 10.
5. Add a new column 'Total' that calculates Price \* Quantity for each row.

### 3.4 Solutions

```

1.
1 import pandas as pd
2 import numpy as np
3
4 np.random.seed(42)
5 sales_data = {
6     'Product': [f'Product_{i}' for i in range(1, 11)],
7     'Category': np.random.choice(['Electronics', 'Clothing', 'Books', 'Home
8     ''], 10),
9     'Price': np.random.randint(10, 200, 10),
10    'Quantity': np.random.randint(1, 20, 10),
11    'Date': pd.date_range('2024-01-01', periods=10)
12 }
13 sales_df = pd.DataFrame(sales_data)
14 print(sales_df)

```

```

1 electronics_expensive = sales_df.loc[
2     (sales_df['Category'] == 'Electronics') &
3     (sales_df['Price'] > 50)
4 ]
5 print(electronics_expensive)

```

```

3.
1 first_five = sales_df.iloc[0:5, [0, 2]] # Rows 0-4, columns 0 and 2
2 print(first_five)

```

```

1 low_quantity = sales_df[sales_df['Quantity'] < 10]
2 print(low_quantity)

```

```

4.
1 sales_df['Total'] = sales_df['Price'] * sales_df['Quantity']
2 print(sales_df)

```

## Data Cleaning and Preparation

### 4.1 Handling Missing Data

#### 4.1.1 Detecting Missing Values

```
1 # Create DataFrame with missing values
2 data = {
3     'A': [1, 2, np.nan, 4, 5],
4     'B': [np.nan, 2, 3, np.nan, 5],
5     'C': [1, 2, 3, 4, 5],
6     'D': [np.nan, np.nan, np.nan, np.nan, np.nan]
7 }
8 df = pd.DataFrame(data)
9
10 print("DataFrame with missing values:")
11 print(df)
12
13 print("\nMissing values per column:")
14 print(df.isnull().sum())
15
16 print("\nPercentage of missing values:")
17 print(df.isnull().sum() / len(df) * 100)
18
19 print("\nRows with any missing values:")
20 print(df[df.isnull().any(axis=1)])
21
22 print("\nColumns with any missing values:")
23 print(df.columns[df.isnull().any()].tolist())
```

Listing 4.1: Identifying missing data

#### 4.1.2 Removing Missing Data

```
1 # Remove rows with any missing values
2 df_dropped_rows = df.dropna()
3 print("After dropping rows with any NaN:")
4 print(df_dropped_rows)
5
6 # Remove columns with any missing values
7 df_dropped_cols = df.dropna(axis=1)
8 print("\nAfter dropping columns with any NaN:")
9 print(df_dropped_cols)
10
11 # Remove rows where all values are NaN
12 df_dropped_all_na = df.dropna(how='all')
13 print("\nAfter dropping rows where all values are NaN:")
14 print(df_dropped_all_na)
15
16 # Remove rows with less than 3 non-NaN values
```

```
17 df_dropped_thresh = df.dropna(thresh=3)
18 print("\nAfter dropping rows with less than 3 non-NaN values:")
19 print(df_dropped_thresh)
```

Listing 4.2: Removing missing values

### 4.1.3 Filling Missing Data

```
1 # Fill with specific value
2 df_filled_0 = df.fillna(0)
3 print("Filled with 0:")
4 print(df_filled_0)
5
6 # Fill with mean of each column
7 df_filled_mean = df.fillna(df.mean())
8 print("\nFilled with column means:")
9 print(df_filled_mean)
10
11 # Forward fill
12 df_ffill = df.fillna(method='ffill')
13 print("\nForward filled:")
14 print(df_ffill)
15
16 # Backward fill
17 df_bfill = df.fillna(method='bfill')
18 print("\nBackward filled:")
19 print(df_bfill)
20
21 # Different fill values for different columns
22 fill_values = {'A': 0, 'B': df['B'].mean(), 'C': 999}
23 df_custom_fill = df.fillna(fill_values)
24 print("\nCustom filled:")
25 print(df_custom_fill)
26
27 # Using interpolate for time series
28 df_interpolated = df.interpolate()
29 print("\nInterpolated:")
30 print(df_interpolated)
```

Listing 4.3: Filling missing values

## 4.2 Data Type Conversion

```
1 # Create sample data with mixed types
2 data = {
3     'integers_as_strings': ['1', '2', '3', '4', '5'],
4     'floats_as_strings': ['1.1', '2.2', '3.3', '4.4', '5.5'],
5     'booleans_as_strings': ['True', 'False', 'True', 'False', 'True'],
6     'dates_as_strings': ['2024-01-01', '2024-01-02', '2024-01-03', '2024-01-04',
7         '2024-01-05'],
8     'mixed_numeric': [1, '2', 3.0, '4', 5]
9 }
10 df_types = pd.DataFrame(data)
11
12 print("Original data types:")
13 print(df_types.dtypes)
```

```

13
14 # Convert to appropriate types
15 df_types['integers_as_strings'] = df_types['integers_as_strings'].astype(int)
16 df_types['floats_as_strings'] = df_types['floats_as_strings'].astype(float)
17 df_types['booleans_as_strings'] = df_types['booleans_as_strings'].map({'True':
    True, 'False': False})
18 df_types['dates_as_strings'] = pd.to_datetime(df_types['dates_as_strings'])
19 df_types['mixed_numeric'] = pd.to_numeric(df_types['mixed_numeric'], errors='
    coerce')
20
21 print("\nAfter conversion:")
22 print(df_types.dtypes)
23 print(df_types)

```

Listing 4.4: Converting data types

## 4.3 String Operations

```

1 # Create sample text data
2 text_data = {
3     'name': ['john doe', 'JANE SMITH', 'Bob Johnson', 'alice brown', 'CHARLIE
    WILSON'],
4     'email': ['john@email.com', 'jane@test.org', 'bob@company.com', 'alice@home.
    net', 'charlie@work.io'],
5     'description': ['Python developer with 5 years experience',
6                     'Data scientist specializing in ML',
7                     'Senior software engineer',
8                     'Junior web developer',
9                     'Project manager with technical background']
10 }
11 df_text = pd.DataFrame(text_data)
12
13 # String methods
14 df_text['name_clean'] = df_text['name'].str.title()
15 df_text['email_domain'] = df_text['email'].str.split('@').str[1]
16 df_text['name_length'] = df_text['name'].str.len()
17 df_text['has_experience'] = df_text['description'].str.contains('experience',
    case=False)
18 df_text['words_count'] = df_text['description'].str.split().str.len()
19
20 print("Text processing results:")
21 print(df_text)

```

Listing 4.5: String manipulation with pandas

## 4.4 Exercises

1. Create a DataFrame with various types of missing data and practice different handling techniques.
2. Convert a column of string dates to datetime objects and extract year, month, and day.
3. Clean a column of messy string data (mixed case, extra spaces) and standardize the format.
4. Handle a column with mixed numeric and string values, converting everything to numeric.

5. Create a function to automatically detect and report data quality issues in a DataFrame.

## 4.5 Solutions

```
1.
1 # Create DataFrame with missing data
2 missing_data = {
3     'col1': [1, np.nan, 3, None, 5],
4     'col2': ['a', 'b', None, 'd', np.nan],
5     'col3': [1.1, 2.2, np.nan, np.nan, 5.5]
6 }
7 df_missing = pd.DataFrame(missing_data)
8
9 # Different handling methods
10 df_dropped = df_missing.dropna()
11 df_filled = df_missing.fillna({'col1': 0, 'col2': 'missing', 'col3':
12     df_missing['col3'].mean()})
12 df_ffilled = df_missing.ffill()
```

```
1 dates = pd.Series(['2024-01-15', '2024-02-20', '2024-03-25', '2024-04-30'])
2 date_series = pd.to_datetime(dates)
3 result = pd.DataFrame({
4     'original': dates,
5     'datetime': date_series,
6     'year': date_series.dt.year,
7     'month': date_series.dt.month,
8     'day': date_series.dt.day
9 })
10 print(result)
```

```
3.
1 messy_data = pd.Series([' JOHN doe ', 'jane SMITH ', ' bob ', 'ALICE',
2     ' Charlie Wilson '])
3 cleaned = messy_data.str.strip().str.title()
4 print(cleaned)
```

```
1 mixed_data = pd.Series([1, '2', 3.0, '4', 'five', 6])
2 numeric_data = pd.to_numeric(mixed_data, errors='coerce')
3 print(numeric_data)
```

```
4.
1 def data_quality_report(df):
2     report = {
3         'missing_values': df.isnull().sum(),
4         'missing_percentage': (df.isnull().sum() / len(df)) * 100,
5         'data_types': df.dtypes,
6         'unique_values': df.nunique(),
7         'memory_usage': df.memory_usage(deep=True)
8     }
9     return pd.DataFrame(report)
10
11 quality_report = data_quality_report(df_missing)
12 print(quality_report)
```

## Part II

# Intermediate Pandas: Data Manipulation

Anshuman Singh — anshuman@datasciencecampus.com



## Data Input/Output Operations

### 5.1 Reading Different File Formats

Pandas provides excellent support for reading data from various file formats. This section covers the most common ones.

```
1 import pandas as pd
2
3 # Reading CSV files (most common)
4 df_csv = pd.read_csv('data.csv')
5 df_csv_custom = pd.read_csv('data.csv', sep=',', header=0, index_col=0)
6
7 # Reading Excel files
8 df_excel = pd.read_excel('data.xlsx', sheet_name='Sheet1')
9 df_excel_multiple = pd.read_excel('data.xlsx', sheet_name=['Sheet1', 'Sheet2'])
10
11 # Reading JSON files
12 df_json = pd.read_json('data.json')
13 df_json_lines = pd.read_json('data.jsonl', lines=True)
14
15 # Reading from SQL database
16 import sqlite3
17 conn = sqlite3.connect('database.db')
18 df_sql = pd.read_sql_query('SELECT * FROM table_name', conn)
19
20 # Reading from URL
21 df_url = pd.read_csv('https://example.com/data.csv')
22
23 # Reading with specific encoding
24 df_encoded = pd.read_csv('data.csv', encoding='utf-8')
```

Listing 5.1: Reading various file formats with pandas

### 5.2 Writing Data to Files

```
1 # Writing to CSV
2 df.to_csv('output.csv', index=False)
3
4 # Writing to Excel
5 df.to_excel('output.xlsx', sheet_name='Data', index=False)
6
7 # Writing to JSON
8 df.to_json('output.json', orient='records')
9
10 # Writing to SQL
11 from sqlalchemy import create_engine
12 engine = create_engine('sqlite:///database.db')
13 df.to_sql('table_name', engine, if_exists='replace', index=False)
14
```

```

15 # Writing to Parquet (efficient for large datasets)
16 df.to_parquet('output.parquet')
17
18 # Writing with compression
19 df.to_csv('output.csv.gz', compression='gzip')

```

Listing 5.2: Writing DataFrames to various file formats

## 5.3 File Reading Parameters and Options

```

1 # CSV reading with various options
2 df = pd.read_csv('data.csv',
3                 sep=',', # Separator
4                 header=0, # Row to use as column names
5                 index_col=0, # Column to use as index
6                 usecols=[0, 1, 2], # Columns to read
7                 nrows=1000, # Number of rows to read
8                 skiprows=[0, 2], # Rows to skip
9                 na_values=['NA', 'null', ''], # Values to treat as NaN
10                dtype={'column1': str, 'column2': int}, # Data types
11                parse_dates=['date_column'], # Parse dates
12                encoding='utf-8' # File encoding
13            )
14
15 # Excel reading options
16 df = pd.read_excel('data.xlsx',
17                   sheet_name=0, # Sheet name or index
18                   header=0, # Row for column names
19                   index_col=0, # Column for index
20                   usecols='A:C,E', # Columns to read
21                   na_values=[''], # Values to treat as NaN
22                   dtype={'col1': str} # Column data types
23            )

```

Listing 5.3: Common parameters for file reading

## 5.4 Handling Large Files

```

1 # Reading in chunks
2 chunk_size = 10000
3 chunks = []
4 for chunk in pd.read_csv('large_file.csv', chunksize=chunk_size):
5     # Process each chunk
6     processed_chunk = chunk[chunk['important_column'] > threshold]
7     chunks.append(processed_chunk)
8
9 # Combine processed chunks
10 result = pd.concat(chunks, ignore_index=True)
11
12 # Using dtypes for memory optimization
13 dtype_dict = {
14     'id': 'int32',
15     'category': 'category',
16     'value': 'float32',
17     'description': 'string'
18 }

```

```

19 df_optimized = pd.read_csv('data.csv', dtype=dtype_dict)
20
21 # Specifying columns to read
22 columns_to_use = ['id', 'name', 'value']
23 df_subset = pd.read_csv('data.csv', usecols=columns_to_use)

```

Listing 5.4: Strategies for large datasets

## 5.5 Exercises

1. Download a sample CSV file from a public URL and read it into pandas.
2. Read an Excel file with multiple sheets and combine them into one DataFrame.
3. Write a DataFrame to both CSV and JSON formats with different formatting options.
4. Create a function that reads a large CSV file in chunks and performs basic analysis on each chunk.
5. Compare the file sizes and read times for CSV, Parquet, and compressed CSV formats.

## 5.6 Solutions

```

1.
1 import pandas as pd
2
3 # Read from URL
4 url = 'https://raw.githubusercontent.com/datasets/iris/master/data/iris.csv'
5 df_iris = pd.read_csv(url)
6 print(f"Loaded {len(df_iris)} rows")
7 print(df_iris.head())

```

```

1 # Read multiple sheets and combine
2 excel_file = 'data.xlsx'
3 sheets_dict = pd.read_excel(excel_file, sheet_name=None)
4
5 # Combine all sheets
6 all_sheets = []
7 for sheet_name, sheet_data in sheets_dict.items():
8     sheet_data['sheet_name'] = sheet_name
9     all_sheets.append(sheet_data)
10
11 combined_df = pd.concat(all_sheets, ignore_index=True)
12 print(combined_df)

```

```

3.
1 # Sample DataFrame
2 data = {'Name': ['Alice', 'Bob'], 'Age': [25, 30], 'City': ['NY', 'LA']}
3 df = pd.DataFrame(data)
4
5 # Write to CSV with custom options
6 df.to_csv('output.csv', index=False, sep='|')
7
8 # Write to JSON with different orientations
9 df.to_json('output_records.json', orient='records', indent=2)
10 df.to_json('output_values.json', orient='values')

```

```

1 def analyze_large_file(file_path, chunk_size=1000):
2     chunk_stats = []
3
4     for i, chunk in enumerate(pd.read_csv(file_path, chunksize=chunk_size)):
5         :
6         stats = {
7             'chunk': i + 1,
8             'rows': len(chunk),
9             'columns': len(chunk.columns),
10            'memory_mb': chunk.memory_usage(deep=True).sum() / 1024**2
11        }
12        chunk_stats.append(stats)
13        print(f"Processed chunk {i+1}: {len(chunk)} rows")
14
15    return pd.DataFrame(chunk_stats)
16
17 # Usage
18 stats_df = analyze_large_file('large_data.csv')
19 print(stats_df)

```

4.

```

1 import time
2 import os
3
4 # Create sample data
5 large_df = pd.DataFrame(np.random.randn(100000, 10),
6                          columns=[f'col_{i}' for i in range(10)])
7
8 # Test different formats
9 formats = ['csv', 'parquet', 'csv.gz']
10 results = []
11
12 for fmt in formats:
13     filename = f'test_data.{fmt}'
14
15     # Write
16     start_write = time.time()
17     if fmt == 'csv':
18         large_df.to_csv(filename, index=False)
19     elif fmt == 'parquet':
20         large_df.to_parquet(filename)
21     elif fmt == 'csv.gz':
22         large_df.to_csv(filename, index=False, compression='gzip')
23     write_time = time.time() - start_write
24
25     # Read
26     start_read = time.time()
27     if fmt == 'csv':
28         pd.read_csv(filename)
29     elif fmt == 'parquet':
30         pd.read_parquet(filename)
31     elif fmt == 'csv.gz':
32         pd.read_csv(filename, compression='gzip')
33     read_time = time.time() - start_read
34
35     # File size
36     file_size = os.path.getsize(filename) / 1024 # KB
37
38     results.append({
39         'format': fmt,

```

```
40         'file_size_kb': file_size,
41         'write_time': write_time,
42         'read_time': read_time
43     })
44
45     # Clean up
46     os.remove(filename)
47
48 results_df = pd.DataFrame(results)
49 print(results_df)
```

## Data Aggregation and Group Operations

### 6.1 GroupBy Fundamentals

```
1 # Sample sales data
2 np.random.seed(42)
3 sales_data = {
4     'Region': np.random.choice(['North', 'South', 'East', 'West'], 100),
5     'Product': np.random.choice(['Widget A', 'Widget B', 'Gadget C'], 100),
6     'Salesperson': np.random.choice(['Alice', 'Bob', 'Charlie', 'Diana'], 100),
7     'Sales': np.random.randint(100, 5000, 100),
8     'Quantity': np.random.randint(1, 50, 100),
9     'Date': pd.date_range('2024-01-01', periods=100, freq='D')
10 }
11 sales_df = pd.DataFrame(sales_data)
12
13 # Basic groupby
14 print("Total sales by region:")
15 region_sales = sales_df.groupby('Region')['Sales'].sum()
16 print(region_sales)
17
18 print("\nAverage sales by product:")
19 product_avg = sales_df.groupby('Product')['Sales'].mean()
20 print(product_avg)
21
22 # Multiple columns grouping
23 print("\nSales by region and product:")
24 region_product = sales_df.groupby(['Region', 'Product'])['Sales'].sum()
25 print(region_product)
```

Listing 6.1: Basic GroupBy operations

### 6.2 Aggregation Methods

```
1 # Multiple aggregation functions
2 agg_results = sales_df.groupby('Region').agg({
3     'Sales': ['sum', 'mean', 'std', 'min', 'max', 'count'],
4     'Quantity': ['sum', 'mean']
5 })
6 print("Multiple aggregations:")
7 print(agg_results)
8
9 # Custom aggregation functions
10 def sales_range(x):
11     return x.max() - x.min()
12
13 def top_sales(x):
```

```

14     return x.nlargest(3).sum()
15
16 custom_agg = sales_df.groupby('Region').agg({
17     'Sales': [sales_range, top_sales, 'mean'],
18     'Quantity': ['sum', lambda x: x.mean()]
19 })
20 print("\nCustom aggregations:")
21 print(custom_agg)
22
23 # Named aggregations (cleaner syntax)
24 named_agg = sales_df.groupby('Region').agg(
25     total_sales=('Sales', 'sum'),
26     avg_sales=('Sales', 'mean'),
27     max_quantity=('Quantity', 'max'),
28     unique_products=('Product', 'nunique')
29 )
30 print("\nNamed aggregations:")
31 print(named_agg)

```

Listing 6.2: Advanced aggregation techniques

## 6.3 Transform and Filter Operations

```

1 # Transform - apply function to each group and return same shape
2 sales_df['region_avg_sales'] = sales_df.groupby('Region')['Sales'].transform('
    mean')
3 sales_df['sales_vs_region_avg'] = sales_df['Sales'] - sales_df['region_avg_sales
    ']
4
5 print("Data with transformed columns:")
6 print(sales_df[['Region', 'Sales', 'region_avg_sales', 'sales_vs_region_avg']].
    head())
7
8 # Filter - filter groups based on conditions
9 def high_sales_group(group):
10     return group['Sales'].sum() > 50000
11
12 high_sales_regions = sales_df.groupby('Region').filter(high_sales_group)
13 print(f"\nRegions with total sales > 50,000: {high_sales_regions['Region'].
    unique()}")
14
15 # Size and count
16 group_sizes = sales_df.groupby(['Region', 'Product']).size()
17 print("\nGroup sizes:")
18 print(group_sizes)

```

Listing 6.3: Transform and filter in GroupBy

## 6.4 Pivot Tables

```

1 # Basic pivot table
2 pivot_simple = pd.pivot_table(sales_df,
3                               values='Sales',
4                               index='Region',
5                               columns='Product',
6                               aggfunc='sum')

```

```

7 print("Simple pivot table:")
8 print(pivot_simple)
9
10 # Multi-level pivot table
11 pivot_complex = pd.pivot_table(sales_df,
12                                values=['Sales', 'Quantity'],
13                                index=['Region', 'Salesperson'],
14                                columns=['Product'],
15                                aggfunc={'Sales': ['sum', 'mean'], 'Quantity': 'sum'},
16                                fill_value=0,
17                                margins=True)
18 print("\nComplex pivot table:")
19 print(pivot_complex)
20
21 # Using pivot (vs pivot_table)
22 pivot_simple_alt = sales_df.pivot_table(index='Region',
23                                         columns='Product',
24                                         values='Sales',
25                                         aggfunc='mean')
26 print("\nAlternative pivot syntax:")
27 print(pivot_simple_alt)

```

Listing 6.4: Creating pivot tables

## 6.5 Cross Tabulation

```

1 # Create cross tabulation
2 cross_tab = pd.crosstab(sales_df['Region'], sales_df['Product'])
3 print("Cross tabulation - counts:")
4 print(cross_tab)
5
6 # Cross tab with margins and normalization
7 cross_tab_margins = pd.crosstab(sales_df['Region'], sales_df['Product'],
8                                margins=True, margins_name='Total')
9 print("\nCross tabulation with margins:")
10 print(cross_tab_margins)
11
12 # Normalized cross tab
13 cross_tab_norm = pd.crosstab(sales_df['Region'], sales_df['Product'],
14                              normalize='index') # Normalize by row
15 print("\nNormalized cross tabulation (by row):")
16 print(cross_tab_norm)
17
18 # Multi-dimensional cross tab
19 cross_tab_multi = pd.crosstab([sales_df['Region'], sales_df['Salesperson']],
20                               sales_df['Product'])
21 print("\nMulti-dimensional cross tab:")
22 print(cross_tab_multi)

```

Listing 6.5: Cross tabulation for frequency analysis

## 6.6 Exercises

1. Create a DataFrame with student data (name, subject, grade, year) and calculate average grades by subject and year.



2. Use transform to center each student's grades around their personal average.
3. Create a pivot table showing average grades by subject and year.
4. Use crosstab to show the count of students by grade category (A: 90-100, B: 80-89, etc.) for each subject.
5. Write a function that performs multiple aggregations on grouped data and returns a formatted report.

## 6.7 Solutions

```

1.
1 # Student data
2 np.random.seed(42)
3 students = []
4 for year in [2022, 2023, 2024]:
5     for subject in ['Math', 'Science', 'English', 'History']:
6         for _ in range(20):
7             students.append({
8                 'name': f'Student_{np.random.randint(1, 101)}',
9                 'subject': subject,
10                'grade': np.random.randint(60, 101),
11                'year': year
12            })
13
14 students_df = pd.DataFrame(students)
15
16 # Average grades by subject and year
17 avg_grades = students_df.groupby(['subject', 'year'])['grade'].mean().round
18                               (2)
19 print(avg_grades)

```

```

1 # Center grades around student average
2 students_df['student_avg'] = students_df.groupby('name')['grade'].transform
3                               ('mean')
4 students_df['centered_grade'] = students_df['grade'] - students_df['
5     student_avg']
6 print(students_df[['name', 'grade', 'student_avg', 'centered_grade']].head
7       (10))

```

```

3.
1 # Pivot table
2 grades_pivot = pd.pivot_table(students_df,
3                               values='grade',
4                               index='subject',
5                               columns='year',
6                               aggfunc='mean')
7 print(grades_pivot)

```

```

1 # Create grade categories
2 def get_grade_category(grade):
3     if grade >= 90: return 'A'
4     elif grade >= 80: return 'B'
5     elif grade >= 70: return 'C'
6     elif grade >= 60: return 'D'
7     else: return 'F'

```

```
8
9 students_df['grade_category'] = students_df['grade'].apply(
10     get_grade_category)
11 # Cross tabulation
12 grade_crosstab = pd.crosstab(students_df['subject'],
13                               students_df['grade_category'])
14 print(grade_crosstab)
```

```
4.
1 def student_performance_report(df):
2     report = df.groupby(['subject', 'year']).agg({
3         'grade': ['count', 'mean', 'std', 'min', 'max'],
4         'name': 'nunique'
5     }).round(2)
6
7     # Flatten column names
8     report.columns = ['_'.join(col).strip() for col in report.columns.
9                       values]
9     report = report.rename(columns={'name_nunique': 'unique_students'})
10
11     return report
12
13 performance_report = student_performance_report(students_df)
14 print(performance_report)
```

## Data Merging, Joining, and Concatenating

### 7.1 Concatenation Basics

```
1 # Create sample DataFrames
2 df1 = pd.DataFrame({
3     'A': ['A0', 'A1', 'A2', 'A3'],
4     'B': ['B0', 'B1', 'B2', 'B3'],
5     'C': ['C0', 'C1', 'C2', 'C3']
6 }, index=[0, 1, 2, 3])
7
8 df2 = pd.DataFrame({
9     'A': ['A4', 'A5', 'A6', 'A7'],
10    'B': ['B4', 'B5', 'B6', 'B7'],
11    'C': ['C4', 'C5', 'C6', 'C7']
12 }, index=[4, 5, 6, 7])
13
14 df3 = pd.DataFrame({
15     'A': ['A8', 'A9', 'A10', 'A11'],
16     'B': ['B8', 'B9', 'B10', 'B11'],
17     'C': ['C8', 'C9', 'C10', 'C11']
18 }, index=[8, 9, 10, 11])
19
20 # Basic concatenation (vertical)
21 result_vertical = pd.concat([df1, df2, df3])
22 print("Vertical concatenation:")
23 print(result_vertical)
24
25 # Horizontal concatenation
26 df4 = pd.DataFrame({
27     'D': ['D0', 'D1', 'D2', 'D3'],
28     'E': ['E0', 'E1', 'E2', 'E3']
29 }, index=[0, 1, 2, 3])
30
31 result_horizontal = pd.concat([df1, df4], axis=1)
32 print("\nHorizontal concatenation:")
33 print(result_horizontal)
34
35 # Concatenation with keys
36 result_keys = pd.concat([df1, df2, df3], keys=['X', 'Y', 'Z'])
37 print("\nConcatenation with keys:")
38 print(result_keys)
```

Listing 7.1: Concatenating DataFrames

### 7.2 Merging DataFrames

```

1 # Create sample DataFrames for merging
2 employees = pd.DataFrame({
3     'employee_id': [1, 2, 3, 4, 5],
4     'name': ['Alice', 'Bob', 'Charlie', 'Diana', 'Eva'],
5     'department_id': [101, 102, 101, 103, 104]
6 })
7
8 departments = pd.DataFrame({
9     'department_id': [101, 102, 103, 105],
10    'department_name': ['HR', 'Engineering', 'Marketing', 'Sales'],
11    'manager': ['Carol', 'Dave', 'Frank', 'Grace']
12 })
13
14 salaries = pd.DataFrame({
15     'employee_id': [1, 2, 3, 6],
16     'salary': [50000, 60000, 55000, 70000]
17 })
18
19 # Inner join (default)
20 inner_join = pd.merge(employees, departments, on='department_id')
21 print("Inner join:")
22 print(inner_join)
23
24 # Left join
25 left_join = pd.merge(employees, departments, on='department_id', how='left')
26 print("\nLeft join:")
27 print(left_join)
28
29 # Right join
30 right_join = pd.merge(employees, departments, on='department_id', how='right')
31 print("\nRight join:")
32 print(right_join)
33
34 # Outer join
35 outer_join = pd.merge(employees, departments, on='department_id', how='outer')
36 print("\nOuter join:")
37 print(outer_join)
38
39 # Multiple DataFrames
40 multi_merge = pd.merge(employees, departments, on='department_id').merge(
41     salaries, on='employee_id')
42 print("\nMultiple merges:")
43 print(multi_merge)

```

Listing 7.2: Merging DataFrames with different join types

## 7.3 Advanced Merging Techniques

```

1 # Merging on multiple columns
2 df1 = pd.DataFrame({
3     'key1': ['A', 'B', 'C', 'D'],
4     'key2': [1, 2, 3, 4],
5     'value1': ['X1', 'X2', 'X3', 'X4']
6 })
7
8 df2 = pd.DataFrame({
9     'key1': ['A', 'B', 'C', 'E'],

```

```

10     'key2': [1, 2, 3, 5],
11     'value2': ['Y1', 'Y2', 'Y3', 'Y4']
12 })
13
14 merge_multiple_keys = pd.merge(df1, df2, on=['key1', 'key2'])
15 print("Merge on multiple keys:")
16 print(merge_multiple_keys)
17
18 # Merging with different column names
19 df3 = pd.DataFrame({
20     'id': [1, 2, 3, 4],
21     'name': ['Alice', 'Bob', 'Charlie', 'Diana']
22 })
23
24 df4 = pd.DataFrame({
25     'emp_id': [1, 2, 5, 6],
26     'salary': [50000, 60000, 70000, 80000]
27 })
28
29 merge_different_names = pd.merge(df3, df4, left_on='id', right_on='emp_id')
30 print("\nMerge with different column names:")
31 print(merge_different_names)
32
33 # Indicator for merge source
34 merge_with_indicator = pd.merge(employees, departments, on='department_id',
35                                 how='outer', indicator=True)
36 print("\nMerge with indicator:")
37 print(merge_with_indicator)
38
39 # Suffixes for overlapping columns
40 df5 = pd.DataFrame({
41     'key': [1, 2, 3],
42     'value': [100, 200, 300]
43 })
44
45 df6 = pd.DataFrame({
46     'key': [1, 2, 4],
47     'value': [400, 500, 600]
48 })
49
50 merge_suffixes = pd.merge(df5, df6, on='key', suffixes=('_left', '_right'))
51 print("\nMerge with suffixes:")
52 print(merge_suffixes)

```

Listing 7.3: Advanced merging scenarios

## 7.4 Joining DataFrames

```

1 # Create DataFrames with meaningful indexes
2 left_df = pd.DataFrame({
3     'A': ['A0', 'A1', 'A2'],
4     'B': ['B0', 'B1', 'B2']
5 }, index=['K0', 'K1', 'K2'])
6
7 right_df = pd.DataFrame({
8     'C': ['C0', 'C2', 'C3'],
9     'D': ['D0', 'D2', 'D3']
10 }, index=['K0', 'K2', 'K3'])

```

```

11
12 # Left join on index
13 left_join_index = left_df.join(right_df, how='left')
14 print("Left join on index:")
15 print(left_join_index)
16
17 # Right join on index
18 right_join_index = left_df.join(right_df, how='right')
19 print("\nRight join on index:")
20 print(right_join_index)
21
22 # Inner join on index
23 inner_join_index = left_df.join(right_df, how='inner')
24 print("\nInner join on index:")
25 print(inner_join_index)
26
27 # Outer join on index
28 outer_join_index = left_df.join(right_df, how='outer')
29 print("\nOuter join on index:")
30 print(outer_join_index)
31
32 # Join with multiple DataFrames
33 another_df = pd.DataFrame({
34     'E': ['E0', 'E1', 'E2']
35 }, index=['K0', 'K1', 'K3'])
36
37 multi_join = left_df.join([right_df, another_df])
38 print("\nJoin with multiple DataFrames:")
39 print(multi_join)

```

Listing 7.4: Joining DataFrames on index

## 7.5 Exercises

1. Create three DataFrames with customer, order, and product data and merge them appropriately.
2. Concatenate multiple monthly sales reports into a single DataFrame.
3. Perform a merge that shows which customers have no orders (anti-join).
4. Create a DataFrame that combines employee information with their department details using different join types.
5. Compare the performance of merge vs join operations on large datasets.

## 7.6 Solutions

```

1.
1 # Sample data
2 customers = pd.DataFrame({
3     'customer_id': [1, 2, 3, 4],
4     'name': ['Alice', 'Bob', 'Charlie', 'Diana'],
5     'city': ['NY', 'LA', 'Chicago', 'Miami']
6 })
7
8 orders = pd.DataFrame({

```

```

9     'order_id': [101, 102, 103, 104],
10     'customer_id': [1, 2, 1, 5],
11     'product_id': [201, 202, 203, 201],
12     'amount': [100, 200, 150, 300]
13 })
14
15 products = pd.DataFrame({
16     'product_id': [201, 202, 203, 204],
17     'product_name': ['Widget', 'Gadget', 'Tool', 'Accessory'],
18     'price': [50, 100, 75, 25]
19 })
20
21 # Merge all three
22 result = pd.merge(customers, orders, on='customer_id', how='inner')
23 result = pd.merge(result, products, on='product_id', how='left')
24 print(result)

```

```

1 # Create monthly reports
2 jan_sales = pd.DataFrame({
3     'product': ['A', 'B', 'C'],
4     'jan_sales': [100, 150, 200]
5 })
6
7 feb_sales = pd.DataFrame({
8     'product': ['A', 'B', 'D'],
9     'feb_sales': [120, 160, 180]
10 })
11
12 mar_sales = pd.DataFrame({
13     'product': ['A', 'C', 'D'],
14     'mar_sales': [110, 210, 190]
15 })
16
17 # Concatenate vertically
18 all_sales = pd.concat([jan_sales, feb_sales, mar_sales], ignore_index=True)
19 print("Vertical concatenation:")
20 print(all_sales)
21
22 # Merge horizontally
23 quarterly_sales = jan_sales.merge(feb_sales, on='product', how='outer')
24 quarterly_sales = quarterly_sales.merge(mar_sales, on='product', how='outer')
25 print("\nHorizontal merge:")
26 print(quarterly_sales)

```

3.

```

1 # Anti-join: customers with no orders
2 customers_with_orders = pd.merge(customers, orders, on='customer_id', how='inner')
3 all_customers = pd.merge(customers, orders, on='customer_id', how='left',
4     indicator=True)
5 customers_no_orders = all_customers[all_customers['_merge'] == 'left_only']
6 print("Customers with no orders:")
7 print(customers_no_orders[['customer_id', 'name']])

```

```

1 employees = pd.DataFrame({
2     'emp_id': [1, 2, 3, 4, 5],
3     'name': ['Alice', 'Bob', 'Charlie', 'Diana', 'Eva'],
4     'dept_id': [101, 102, 101, 103, 104]

```

```

5 })
6
7 departments = pd.DataFrame({
8     'dept_id': [101, 102, 103, 105],
9     'dept_name': ['HR', 'Engineering', 'Marketing', 'Sales']
10 })
11
12 # Different join types
13 inner_result = pd.merge(employees, departments, on='dept_id', how='inner')
14 left_result = pd.merge(employees, departments, on='dept_id', how='left')
15 right_result = pd.merge(employees, departments, on='dept_id', how='right')
16 outer_result = pd.merge(employees, departments, on='dept_id', how='outer')
17
18 print("Inner join result:")
19 print(inner_result)
20 print(f"Rows: {len(inner_result)}")

```

```

4.
1 import time
2
3 # Create larger datasets for performance comparison
4 large_df1 = pd.DataFrame({
5     'key': range(10000),
6     'value1': np.random.randn(10000)
7 })
8
9 large_df2 = pd.DataFrame({
10     'key': range(5000, 15000),
11     'value2': np.random.randn(10000)
12 })
13
14 # Time merge operation
15 start_time = time.time()
16 merge_result = pd.merge(large_df1, large_df2, on='key', how='inner')
17 merge_time = time.time() - start_time
18
19 # Time join operation (setting index first)
20 large_df1_indexed = large_df1.set_index('key')
21 large_df2_indexed = large_df2.set_index('key')
22
23 start_time = time.time()
24 join_result = large_df1_indexed.join(large_df2_indexed, how='inner')
25 join_time = time.time() - start_time
26
27 print(f"Merge time: {merge_time:.4f} seconds")
28 print(f"Join time: {join_time:.4f} seconds")
29 print(f"Merge result shape: {merge_result.shape}")
30 print(f"Join result shape: {join_result.shape}")

```



## Part III

# Advanced Pandas: Optimization and Real-World Applications

Anshuman Singh — anshuman.singh@datacamp.com

## Advanced Data Manipulation with MultiIndex

### 8.1 Creating and Understanding MultiIndex

```
1 # Creating MultiIndex from arrays
2 arrays = [
3     ['Bar', 'Bar', 'Baz', 'Baz', 'Foo', 'Foo', 'Qux', 'Qux'],
4     ['one', 'two', 'one', 'two', 'one', 'two', 'one', 'two']
5 ]
6 tuples = list(zip(*arrays))
7 index = pd.MultiIndex.from_tuples(tuples, names=['first', 'second'])
8 df_multi = pd.DataFrame(np.random.randn(8, 4), index=index, columns=['A', 'B', 'C', 'D'])
9
10 print("MultiIndex DataFrame:")
11 print(df_multi)
12 print(f"\nIndex names: {df_multi.index.names}")
13 print(f"Index levels: {df_multi.index.levels}")
14
15 # Creating MultiIndex from product
16 index_product = pd.MultiIndex.from_product(
17     [['Group1', 'Group2'], ['SubgroupA', 'SubgroupB'], ['Item1', 'Item2']],
18     names=['Group', 'Subgroup', 'Item']
19 )
20 df_product = pd.DataFrame(np.random.randn(8, 3), index=index_product, columns=['X', 'Y', 'Z'])
21 print("\nMultiIndex from product:")
22 print(df_product)
```

Listing 8.1: Working with MultiIndex DataFrames

### 8.2 Indexing with MultiIndex

```
1 # Different ways to access MultiIndex data
2 print("First level indexing:")
3 print(df_multi.loc['Bar'])
4
5 print("\nMultiple first level indices:")
6 print(df_multi.loc[['Bar', 'Foo']])
7
8 print("\nSpecific multi-level index:")
9 print(df_multi.loc[('Bar', 'one')])
10
11 print("\nSlicing with MultiIndex:")
12 print(df_multi.loc[('Bar', 'one'):(('Baz', 'two'))])
13
14 # Using xs for cross-section
```

```

15 print("\nCross-section at first level:")
16 print(df_multi.xs('Bar', level='first'))
17
18 print("\nCross-section at second level:")
19 print(df_multi.xs('one', level='second'))
20
21 # Boolean indexing with MultiIndex
22 print("\nBoolean indexing:")
23 print(df_multi[df_multi['A'] > 0])

```

Listing 8.2: Advanced MultiIndex operations

## 8.3 Reshaping with Stack and Unstack

```

1 # Sample data for reshaping
2 data = {
3     'Year': [2022, 2022, 2023, 2023, 2022, 2022, 2023, 2023],
4     'Quarter': ['Q1', 'Q2', 'Q1', 'Q2', 'Q1', 'Q2', 'Q1', 'Q2'],
5     'Region': ['North', 'North', 'North', 'North', 'South', 'South', 'South', 'South'],
6     'Sales': [100, 150, 120, 180, 90, 130, 110, 160],
7     'Profit': [20, 30, 25, 35, 18, 26, 22, 32]
8 }
9 df_sales = pd.DataFrame(data)
10
11 # Pivot to create MultiIndex columns
12 pivoted = df_sales.pivot_table(
13     index=['Year', 'Quarter'],
14     columns='Region',
15     values=['Sales', 'Profit']
16 )
17 print("Pivoted DataFrame:")
18 print(pivoted)
19
20 # Stack operation
21 stacked = pivoted.stack()
22 print("\nStacked DataFrame:")
23 print(stacked)
24
25 # Unstack operation
26 unstacked = stacked.unstack()
27 print("\nUnstacked DataFrame:")
28 print(unstacked)
29
30 # Unstack specific level
31 unstack_region = stacked.unstack('Region')
32 print("\nUnstacked by Region:")
33 print(unstack_region)

```

Listing 8.3: Stack and unstack operations

## 8.4 Advanced MultiIndex Operations

```

1 # Setting and resetting index
2 df_reset = df_multi.reset_index()
3 print("After reset_index:")

```

```

4 print(df_reset)
5
6 df_set = df_reset.set_index(['first', 'second'])
7 print("\nAfter set_index:")
8 print(df_set)
9
10 # Swapping index levels
11 df_swapped = df_multi.swaplevel('first', 'second')
12 print("\nAfter swapping levels:")
13 print(df_swapped)
14
15 # Sorting index
16 df_sorted = df_multi.sort_index(level='first')
17 print("\nAfter sorting by first level:")
18 print(df_sorted)
19
20 # MultiIndex on columns
21 df_multi_cols = pd.DataFrame(
22     np.random.randn(4, 8),
23     columns=pd.MultiIndex.from_product(['Group1', 'Group2'], ['A', 'B', 'C', 'D'])
24 )
25 print("\nMultiIndex columns:")
26 print(df_multi_cols)

```

Listing 8.4: Complex MultiIndex manipulations

## 8.5 Exercises

1. Create a MultiIndex DataFrame representing sales data across different regions, products, and time periods.
2. Use stack/unstack to reshape the data between wide and long formats.
3. Calculate summary statistics for each level of the MultiIndex.
4. Filter the MultiIndex DataFrame to show only specific combinations of index levels.
5. Compare the performance of operations on MultiIndex vs single-index DataFrames.

## 8.6 Solutions

```

1.
1 # Create comprehensive sales data
2 index = pd.MultiIndex.from_product([
3     ['North', 'South', 'East', 'West'],
4     ['Electronics', 'Clothing', 'Home'],
5     ['Q1', 'Q2', 'Q3', 'Q4']
6 ], names=['Region', 'Category', 'Quarter'])
7
8 sales_data = pd.DataFrame({
9     'Revenue': np.random.randint(10000, 100000, 48),
10    'Units': np.random.randint(100, 1000, 48),
11    'Profit': np.random.randint(1000, 20000, 48)
12 }, index=index)
13
14 print(sales_data.head(10))

```

```

1 # Reshape using stack/unstack
2 # Wide to long
3 long_format = sales_data.stack()
4 print("Long format (stacked):")
5 print(long_format.head(10))
6
7 # Long to wide
8 wide_format = long_format.unstack()
9 print("\nWide format (unstacked):")
10 print(wide_format.head())

```

3.

```

1 # Summary statistics by different levels
2 print("By Region:")
3 print(sales_data.groupby(level='Region').mean())
4
5 print("\nBy Category:")
6 print(sales_data.groupby(level='Category').sum())
7
8 print("\nBy Region and Category:")
9 print(sales_data.groupby(level=['Region', 'Category']).agg(['mean', 'sum',
    'std']))

```

```

1 # Filter specific combinations
2 # Using loc with tuples
3 filtered = sales_data.loc[
4     (['North', 'South'], ['Electronics', 'Clothing'], ['Q1', 'Q2'])
5 ]
6 print("Filtered data:")
7 print(filtered)
8
9 # Using query (if index levels are properly named)
10 # filtered_query = sales_data.query('Region in ["North", "South"] and
    Quarter in ["Q1", "Q2"]')

```

4.

```

1 import time
2
3 # Create equivalent single-index DataFrame
4 sales_data_single = sales_data.reset_index()
5
6 # Compare grouping operations
7 start_time = time.time()
8 multi_result = sales_data.groupby(level=['Region', 'Category']).mean()
9 multi_time = time.time() - start_time
10
11 start_time = time.time()
12 single_result = sales_data_single.groupby(['Region', 'Category']).mean()
13 single_time = time.time() - start_time
14
15 print(f"MultiIndex grouping time: {multi_time:.6f} seconds")
16 print(f"Single-index grouping time: {single_time:.6f} seconds")
17 print(f"MultiIndex is {single_time/multi_time:.2f}x faster")

```

## Performance Optimization and Vectorization

### 9.1 Understanding Vectorization

```
1 import pandas as pd
2 import numpy as np
3 import time
4
5 # Create large dataset
6 np.random.seed(42)
7 large_df = pd.DataFrame({
8     'A': np.random.rand(100000),
9     'B': np.random.rand(100000),
10    'C': np.random.rand(100000)
11 })
12
13 # Method 1: Using apply (slow)
14 def slow_calculation(row):
15     return row['A'] * row['B'] + row['C']
16
17 start_time = time.time()
18 large_df['slow_result'] = large_df.apply(slow_calculation, axis=1)
19 slow_time = time.time() - start_time
20
21 # Method 2: Vectorized operations (fast)
22 start_time = time.time()
23 large_df['fast_result'] = large_df['A'] * large_df['B'] + large_df['C']
24 fast_time = time.time() - start_time
25
26 print(f"Apply method: {slow_time:.4f} seconds")
27 print(f"Vectorized method: {fast_time:.4f} seconds")
28 print(f"Speedup: {slow_time/fast_time:.2f}x")
29
30 # Verify results are the same
31 print(f"Results match: {np.allclose(large_df['slow_result'], large_df['fast_result'])}")
```

Listing 9.1: Vectorized operations vs loops

### 9.2 Memory Optimization

```
1 # Check memory usage
2 print("Original memory usage:")
3 print(large_df.memory_usage(deep=True))
4 print(f"Total: {large_df.memory_usage(deep=True).sum() / 1024**2:.2f} MB")
5
6 # Optimize data types
```

```

7 def optimize_dtypes(df):
8     optimized_df = df.copy()
9
10    # Integer columns
11    int_cols = optimized_df.select_dtypes(include=['int']).columns
12    for col in int_cols:
13        c_min = optimized_df[col].min()
14        c_max = optimized_df[col].max()
15
16        if c_min > np.iinfo(np.int8).min and c_max < np.iinfo(np.int8).max:
17            optimized_df[col] = optimized_df[col].astype(np.int8)
18        elif c_min > np.iinfo(np.int16).min and c_max < np.iinfo(np.int16).max:
19            optimized_df[col] = optimized_df[col].astype(np.int16)
20        elif c_min > np.iinfo(np.int32).min and c_max < np.iinfo(np.int32).max:
21            optimized_df[col] = optimized_df[col].astype(np.int32)
22        else:
23            optimized_df[col] = optimized_df[col].astype(np.int64)
24
25    # Float columns
26    float_cols = optimized_df.select_dtypes(include=['float']).columns
27    for col in float_cols:
28        optimized_df[col] = optimized_df[col].astype(np.float32)
29
30    # Object columns to category if low cardinality
31    obj_cols = optimized_df.select_dtypes(include=['object']).columns
32    for col in obj_cols:
33        num_unique = optimized_df[col].nunique()
34        num_total = len(optimized_df[col])
35        if num_unique / num_total < 0.5: # If less than 50% unique values
36            optimized_df[col] = optimized_df[col].astype('category')
37
38    return optimized_df
39
40    # Create DataFrame with mixed types for demonstration
41    mixed_df = pd.DataFrame({
42        'small_ints': np.random.randint(1, 100, 10000),
43        'large_ints': np.random.randint(1, 1000000, 10000),
44        'floats': np.random.randn(10000),
45        'categories': np.random.choice(['A', 'B', 'C', 'D'], 10000),
46        'strings': np.random.choice(['cat', 'dog', 'bird', 'fish'], 10000)
47    })
48
49    print("\nBefore optimization:")
50    print(mixed_df.memory_usage(deep=True))
51    print(f"Total: {mixed_df.memory_usage(deep=True).sum() / 1024**2:.2f} MB")
52
53    optimized_df = optimize_dtypes(mixed_df)
54
55    print("\nAfter optimization:")
56    print(optimized_df.memory_usage(deep=True))
57    print(f"Total: {optimized_df.memory_usage(deep=True).sum() / 1024**2:.2f} MB")
58    print(f"Memory reduction: {(1 - optimized_df.memory_usage(deep=True).sum() /
59        mixed_df.memory_usage(deep=True).sum()) * 100:.1f}%")

```

Listing 9.2: Reducing memory usage

## 9.3 Efficient Data Processing

```

1 # Method chaining for efficiency
2 # Inefficient: multiple assignments
3 df_inefficient = large_df.copy()
4 df_inefficient = df_inefficient[df_inefficient['A'] > 0.5]
5 df_inefficient['new_col'] = df_inefficient['B'] * 2
6 df_inefficient = df_inefficient.sort_values('C')
7
8 # Efficient: method chaining
9 df_efficient = (large_df[large_df['A'] > 0.5]
10                  .assign(new_col=lambda x: x['B'] * 2)
11                  .sort_values('C'))
12
13 print("Results are identical:", df_inefficient.equals(df_efficient))
14
15 # Using query for complex filtering
16 complex_filtered = large_df.query('A > 0.5 and B < 0.3 and C between 0.1 and 0.9
17                                     ')
18 print(f"Complex filter result shape: {complex_filtered.shape}")
19
20 # Using eval for complex expressions
21 # Note: eval can be faster for very large DataFrames
22 result_eval = large_df.eval('D = A + B * C')
23 print("Using eval:")
24 print(result_eval[['A', 'B', 'C', 'D']].head())

```

Listing 9.3: Efficient data processing techniques

## 9.4 Working with Large Datasets

```

1 # Processing large files in chunks
2 def process_large_file(file_path, chunk_size=10000):
3     """Process a large CSV file in chunks"""
4     processed_chunks = []
5
6     for i, chunk in enumerate(pd.read_csv(file_path, chunksize=chunk_size)):
7         # Process each chunk
8         chunk_processed = (chunk
9                             .query('value > 0') # Filter rows
10                            .assign(processed=True) # Add column
11                            .groupby('category')
12                            .agg({'value': ['sum', 'mean']}) # Aggregate
13                           )
14         processed_chunks.append(chunk_processed)
15
16         if i % 10 == 0: # Progress indicator
17             print(f"Processed {i * chunk_size} rows...")
18
19     # Combine results
20     final_result = pd.concat(processed_chunks)
21     return final_result
22
23 # Using Dask for out-of-core computations (pandas-like interface)
24 try:
25     import dask.dataframe as dd
26
27     # Create a large dataset for demonstration
28     large_csv_path = 'large_dataset.csv'

```



```

29     # Uncomment to create test file:
30     # pd.DataFrame(np.random.randn(1000000, 10)).to_csv(large_csv_path, index=
        False)
31
32     # Read with Dask
33     ddf = dd.read_csv(large_csv_path)
34     print(f"Dask DataFrame: {len(ddf):,} rows")
35
36     # Perform operations (lazy evaluation)
37     result_dask = (ddf[ddf['0'] > 0]
38                   .groupby('1')
39                   .mean()
40                   .compute()) # Actually compute the result
41
42     print("Dask result:")
43     print(result_dask.head())
44
45 except ImportError:
46     print("Dask not installed. Install with: pip install dask")
47
48 # Using pandas with efficient data types from start
49 efficient_dtypes = {
50     'user_id': 'int32',
51     'product_id': 'int32',
52     'category': 'category',
53     'value': 'float32',
54     'timestamp': 'string'
55 }
56
57 # When reading data, specify dtypes
58 # df_efficient = pd.read_csv('data.csv', dtype=efficient_dtypes)

```

Listing 9.4: Strategies for large datasets

## 9.5 Exercises

1. Compare the performance of different methods for calculating row-wise operations on a large DataFrame.
2. Optimize the memory usage of a DataFrame with mixed data types.
3. Process a simulated large dataset using chunk processing and compare with loading the entire dataset.
4. Use method chaining to perform multiple operations in a single expression.
5. Implement a function that automatically suggests optimal data types for a given DataFrame.

## 9.6 Solutions

```

1.
1  # Performance comparison
2  def performance_comparison():
3      df = pd.DataFrame(np.random.rand(100000, 5), columns=list('ABCDE'))
4
5      # Method 1: Iterrows (slowest)
6      start = time.time()

```

```

7     result1 = []
8     for idx, row in df.iterrows():
9         result1.append(row['A'] + row['B'] * row['C'])
10    time1 = time.time() - start
11
12    # Method 2: Apply
13    start = time.time()
14    result2 = df.apply(lambda row: row['A'] + row['B'] * row['C'], axis=1)
15    time2 = time.time() - start
16
17    # Method 3: Vectorized
18    start = time.time()
19    result3 = df['A'] + df['B'] * df['C']
20    time3 = time.time() - start
21
22    print(f"Iterrows: {time1:.4f}s")
23    print(f"Apply: {time2:.4f}s")
24    print(f"Vectorized: {time3:.4f}s")
25    print(f"Vectorized is {time1/time3:.1f}x faster than iterrows")
26
27    performance_comparison()

```

```

1    # Memory optimization
2    def optimize_dataframe(df):
3        # Convert objects to category where appropriate
4        for col in df.select_dtypes(include=['object']):
5            if df[col].nunique() / len(df) < 0.5:
6                df[col] = df[col].astype('category')
7
8        # Downcast numeric columns
9        for col in df.select_dtypes(include=['integer']):
10            df[col] = pd.to_numeric(df[col], downcast='integer')
11
12        for col in df.select_dtypes(include=['float']):
13            df[col] = pd.to_numeric(df[col], downcast='float')
14
15        return df
16
17    # Test
18    test_df = pd.DataFrame({
19        'ints': range(10000),
20        'small_ints': np.random.randint(1, 100, 10000),
21        'floats': np.random.randn(10000),
22        'categories': np.random.choice(['A', 'B', 'C'], 10000),
23        'text': ['hello'] * 10000
24    })
25
26    original_memory = test_df.memory_usage(deep=True).sum()
27    optimized_df = optimize_dataframe(test_df.copy())
28    optimized_memory = optimized_df.memory_usage(deep=True).sum()
29
30    print(f"Original: {original_memory / 1024:.1f} KB")
31    print(f"Optimized: {optimized_memory / 1024:.1f} KB")
32    print(f"Reduction: {(1 - optimized_memory/original_memory)*100:.1f}%")

```

### 3.

```

1    # Chunk processing vs full load
2    def simulate_large_processing():
3        # Create a large CSV for testing

```

```

4     large_data = pd.DataFrame(np.random.randn(500000, 5))
5     large_data.to_csv('test_large.csv', index=False)
6
7     # Method 1: Full load
8     start = time.time()
9     full_df = pd.read_csv('test_large.csv')
10    full_result = full_df[full_df['0'] > 0].mean()
11    full_time = time.time() - start
12
13    # Method 2: Chunk processing
14    start = time.time()
15    chunks = []
16    for chunk in pd.read_csv('test_large.csv', chunksize=10000):
17        filtered_chunk = chunk[chunk['0'] > 0]
18        chunks.append(filtered_chunk)
19    chunk_result = pd.concat(chunks, ignore_index=True).mean()
20    chunk_time = time.time() - start
21
22    print(f"Full load time: {full_time:.2f}s")
23    print(f"Chunk processing time: {chunk_time:.2f}s")
24    print(f"Memory usage full: {full_df.memory_usage().sum() / 1024**2:.1f}
25          MB")
26
27    # Clean up
28    import os
29    os.remove('test_large.csv')
30
31    return full_result, chunk_result
32
33 full_res, chunk_res = simulate_large_processing()
34 print("Results are equal:", full_res.equals(chunk_res))

```

```

1 # Method chaining example
2 result = (pd.DataFrame({
3     'group': np.random.choice(['A', 'B', 'C'], 1000),
4     'value1': np.random.randn(1000),
5     'value2': np.random.randn(1000)
6 })
7     .query('value1 > 0') # Filter positive values
8     .assign( # Create new columns
9         total=lambda x: x['value1'] + x['value2'],
10        ratio=lambda x: x['value1'] / x['value2']
11    )
12    .groupby('group') # Group by category
13    .agg({ # Multiple aggregations
14        'value1': ['mean', 'std'],
15        'value2': ['min', 'max'],
16        'total': 'sum',
17        'ratio': 'mean'
18    })
19    .round(3) # Round results
20    )
21
22 print("Method chaining result:")
23 print(result)

```

4.

```

1 # Data type suggestion function
2 def suggest_optimal_dtypes(df):

```

```

3     suggestions = {}
4
5     for col in df.columns:
6         col_data = df[col]
7         dtype = str(col_data.dtype)
8
9         if col_data.dtype == 'object':
10             unique_ratio = col_data.nunique() / len(col_data)
11             if unique_ratio < 0.5:
12                 suggestions[col] = 'category'
13             else:
14                 suggestions[col] = 'string'
15
16         elif 'int' in dtype:
17             c_min, c_max = col_data.min(), col_data.max()
18             if c_min >= 0: # Unsigned
19                 if c_max < 256:
20                     suggestions[col] = 'uint8'
21                 elif c_max < 65536:
22                     suggestions[col] = 'uint16'
23                 elif c_max < 4294967296:
24                     suggestions[col] = 'uint32'
25             else: # Signed
26                 if c_min > -128 and c_max < 127:
27                     suggestions[col] = 'int8'
28                 elif c_min > -32768 and c_max < 32767:
29                     suggestions[col] = 'int16'
30                 elif c_min > -2147483648 and c_max < 2147483647:
31                     suggestions[col] = 'int32'
32
33         elif 'float' in dtype:
34             suggestions[col] = 'float32'
35
36     return suggestions
37
38 # Test the function
39 sample_df = pd.DataFrame({
40     'small_ints': range(100),
41     'large_ints': range(1000000, 1000100),
42     'floats': np.random.randn(100),
43     'categories': ['A', 'B'] * 50,
44     'strings': [f'text_{i}' for i in range(100)]
45 })
46
47 suggestions = suggest_optimal_dtypes(sample_df)
48 print("Optimal dtype suggestions:")
49 for col, suggested_dtype in suggestions.items():
50     current_dtype = sample_df[col].dtype
51     print(f"{col}: {current_dtype} -> {suggested_dtype}")

```

## Time Series Analysis with Pandas

### 10.1 Working with DateTime Data

```
1 import pandas as pd
2 import numpy as np
3
4 # Creating datetime ranges
5 dates = pd.date_range('2024-01-01', periods=100, freq='D')
6 print("Date range:")
7 print(dates)
8
9 # Creating time series data
10 ts_data = pd.DataFrame({
11     'date': dates,
12     'value': np.random.randn(100).cumsum(),
13     'volume': np.random.randint(1000, 5000, 100)
14 })
15 ts_data = ts_data.set_index('date')
16 print("\nTime series data:")
17 print(ts_data.head())
18
19 # DateTime properties
20 ts_data['year'] = ts_data.index.year
21 ts_data['month'] = ts_data.index.month
22 ts_data['day'] = ts_data.index.day
23 ts_data['dayofweek'] = ts_data.index.dayofweek
24 ts_data['quarter'] = ts_data.index.quarter
25 ts_data['is_weekend'] = ts_data.index.dayofweek.isin([5, 6])
26
27 print("\nData with DateTime properties:")
28 print(ts_data.head())
```

Listing 10.1: DateTime operations in pandas

### 10.2 Time-Based Indexing and Resampling

```
1 # Time-based selection
2 print("January 2024 data:")
3 print(ts_data['2024-01'])
4
5 print("\nFirst week data:")
6 print(ts_data['2024-01-01':'2024-01-07'])
7
8 # Resampling to different frequencies
9 daily_data = ts_data['value']
10 weekly_resampled = daily_data.resample('W').mean()
11 monthly_resampled = daily_data.resample('M').agg(['mean', 'std', 'min', 'max'])
12
13 print("\nWeekly resampled:")
```

```

14 print(weekly_resampled.head())
15
16 print("\nMonthly resampled with multiple aggregations:")
17 print(monthly_resampled.head())
18
19 # Resampling with different methods
20 ohlc_data = daily_data.resample('W').ohlc()
21 print("\nOHLC resampling:")
22 print(ohlc_data.head())

```

Listing 10.2: Time-based operations

## 10.3 Rolling Windows and Expanding Operations

```

1 # Rolling windows
2 ts_data['rolling_mean_7'] = ts_data['value'].rolling(window=7).mean()
3 ts_data['rolling_std_7'] = ts_data['value'].rolling(window=7).std()
4 ts_data['rolling_min_7'] = ts_data['value'].rolling(window=7).min()
5 ts_data['rolling_max_7'] = ts_data['value'].rolling(window=7).max()
6
7 # Expanding windows
8 ts_data['expanding_mean'] = ts_data['value'].expanding().mean()
9 ts_data['expanding_sum'] = ts_data['value'].expanding().sum()
10
11 # Custom rolling functions
12 def custom_roll(x):
13     return (x - x.mean()) / x.std()
14
15 ts_data['z_score_5'] = ts_data['value'].rolling(window=5).apply(custom_roll)
16
17 print("Data with rolling and expanding windows:")
18 print(ts_data[['value', 'rolling_mean_7', 'expanding_mean', 'z_score_5']].head(
    (10)))

```

Listing 10.3: Window operations

## 10.4 Time Series Decomposition

```

1 from statsmodels.tsa.seasonal import seasonal_decompose
2
3 # Create a time series with trend and seasonality
4 np.random.seed(42)
5 t = np.arange(365)
6 trend = 0.1 * t
7 seasonality = 10 * np.sin(2 * np.pi * t / 30)
8 noise = np.random.normal(0, 2, 365)
9 ts_values = trend + seasonality + noise
10
11 # Create DataFrame
12 decomp_df = pd.DataFrame({
13     'value': ts_values,
14     'date': pd.date_range('2024-01-01', periods=365, freq='D')
15 }).set_index('date')
16
17 # Decompose the time series

```

```

18 decomposition = seasonal_decompose(decomp_df['value'], model='additive', period
    =30)
19
20 # Add components to DataFrame
21 decomp_df['trend'] = decomposition.trend
22 decomp_df['seasonal'] = decomposition.seasonal
23 decomp_df['residual'] = decomposition.resid
24
25 print("Decomposition results:")
26 print(decomp_df.head())
27
28 # Check decomposition quality
29 print(f"\nResidual stats - Mean: {decomp_df['residual'].mean():.4f}, Std: {
    decomp_df['residual'].std():.4f}")

```

Listing 10.4: Time series decomposition

## 10.5 Lag Features and Time Shifts

```

1 # Create lag features
2 for lag in [1, 2, 3, 7, 30]:
3     ts_data[f'value_lag_{lag}'] = ts_data['value'].shift(lag)
4
5 # Create forward-looking features (for forecasting)
6 ts_data['value_lead_1'] = ts_data['value'].shift(-1)
7 ts_data['value_lead_7'] = ts_data['value'].shift(-7)
8
9 # Rolling statistics with different windows
10 windows = [3, 7, 14, 30]
11 for window in windows:
12     ts_data[f'rolling_mean_{window}'] = ts_data['value'].rolling(window).mean()
13     ts_data[f'rolling_std_{window}'] = ts_data['value'].rolling(window).std()
14     ts_data[f'rolling_min_{window}'] = ts_data['value'].rolling(window).min()
15     ts_data[f'rolling_max_{window}'] = ts_data['value'].rolling(window).max()
16
17 # Percentage changes
18 ts_data['daily_pct_change'] = ts_data['value'].pct_change()
19 ts_data['weekly_pct_change'] = ts_data['value'].pct_change(7)
20
21 print("Data with lag features and rolling statistics:")
22 print(ts_data.select_dtypes(include=[np.number]).head().round(3))

```

Listing 10.5: Creating lag features

## 10.6 Exercises

1. Create a time series dataset with hourly frequency and add datetime properties.
2. Resample daily data to weekly and monthly frequencies with different aggregation methods.
3. Calculate rolling statistics with multiple window sizes and compare the results.
4. Create lag features for a time series and use them to analyze autocorrelation.
5. Decompose a time series into trend, seasonal, and residual components.

## 10.7 Solutions

1.

```

1 # Hourly time series
2 hourly_dates = pd.date_range('2024-01-01', periods=24*30, freq='H')
3 hourly_data = pd.DataFrame({
4     'timestamp': hourly_dates,
5     'temperature': 20 + 10 * np.sin(2 * np.pi * hourly_dates.hour / 24) +
6         np.random.normal(0, 2, len(hourly_dates)),
7     'humidity': 50 + 20 * np.random.randn(len(hourly_dates))
8 })
9 hourly_data = hourly_data.set_index('timestamp')
10
11 # Add datetime properties
12 hourly_data['hour'] = hourly_data.index.hour
13 hourly_data['day'] = hourly_data.index.day
14 hourly_data['weekday'] = hourly_data.index.day_name()
15 hourly_data['is_night'] = (hourly_data.index.hour < 6) | (hourly_data.index
16     .hour > 22)
17
18 print(hourly_data.head())

```

```

1 # Resample hourly to daily, weekly, monthly
2 daily_resampled = hourly_data.resample('D').agg({
3     'temperature': ['mean', 'min', 'max'],
4     'humidity': 'mean'
5 })
6
7 weekly_resampled = hourly_data.resample('W').agg({
8     'temperature': ['mean', 'std'],
9     'humidity': ['mean', 'std']
10 })
11
12 monthly_resampled = hourly_data.resample('M').agg({
13     'temperature': ['mean', 'min', 'max', 'std'],
14     'humidity': ['mean', 'min', 'max', 'std']
15 })
16
17 print("Daily resampled:")
18 print(daily_resampled.head())

```

2.

```

1 # Multiple rolling windows
2 windows = [6, 12, 24, 168] # 6h, 12h, 24h, 1 week
3 for window in windows:
4     hourly_data[f'temp_rolling_mean_{window}'] = hourly_data['temperature']
5         .rolling(window).mean()
6     hourly_data[f'temp_rolling_std_{window}'] = hourly_data['temperature'].
7         rolling(window).std()
8
9 # Compare smoothing effect
10 comparison = hourly_data[['temperature'] + [f'temp_rolling_mean_{w}' for w
    in windows]].head(200)
11 print("Rolling means comparison:")
12 print(comparison.describe())

```

```

1 # Create lag features and analyze autocorrelation
2 lags = [1, 2, 3, 6, 12, 24] # 1h to 24h lags
3 for lag in lags:

```



```

4     hourly_data[f'temp_lag_{lag}'] = hourly_data['temperature'].shift(lag)
5
6     # Calculate autocorrelation
7     autocorrelations = {}
8     for lag in lags:
9         autocorr = hourly_data['temperature'].autocorr(lag=lag)
10        autocorrelations[lag] = autocorr
11
12    print("Autocorrelations by lag:")
13    for lag, corr in autocorrelations.items():
14        print(f"Lag {lag}h: {corr:.3f}")

```

4.

```

1     # Time series decomposition
2     from statsmodels.tsa.seasonal import seasonal_decompose
3
4     # Use hourly temperature data
5     temp_series = hourly_data['temperature'].iloc[:24*7] # First week
6
7     # Decompose with daily seasonality (24 hours)
8     decomposition = seasonal_decompose(temp_series, model='additive', period
9                                         =24)
10
11    # Create decomposition DataFrame
12    decomp_result = pd.DataFrame({
13        'original': decomposition.observed,
14        'trend': decomposition.trend,
15        'seasonal': decomposition.seasonal,
16        'residual': decomposition.resid
17    })
18
19    print("Decomposition results:")
20    print(decomp_result.head(48)) # First two days
21    print(f"\nResidual statistics:")
22    print(f"Mean: {decomp_result['residual'].mean():.4f}")
23    print(f"Std: {decomp_result['residual'].std():.4f}")

```

## Working with Text Data

### 11.1 Basic String Operations

```
1 # Create sample text data
2 text_data = pd.DataFrame({
3     'name': ['John Doe', 'Jane Smith', 'Bob Johnson', 'Alice Brown', 'Charlie
4             Wilson'],
5     'email': ['john.doe@email.com', 'jane.smith@test.org', 'bob@company.com', '
6               alice.brown@home.net', 'charlie.wilson@work.io'],
7     'description': [
8         'Senior Python developer with 5 years experience in data analysis',
9         'Data scientist specializing in machine learning and AI',
10        'Software engineer with backend development skills',
11        'Junior web developer with frontend experience',
12        'Project manager with technical background and agile methodology'
13    ],
14     'phone': ['(123) 456-7890', '987-654-3210', '555-123-4567', '444-555-6666',
15              '333-222-1111']
16 })
17
18 # Basic string operations
19 text_data['name_upper'] = text_data['name'].str.upper()
20 text_data['name_lower'] = text_data['name'].str.lower()
21 text_data['name_title'] = text_data['name'].str.title()
22 text_data['name_length'] = text_data['name'].str.len()
23 text_data['word_count'] = text_data['description'].str.split().str.len()
24
25 print("\nAfter basic string operations:")
26 print(text_data[['name', 'name_upper', 'name_title', 'name_length', 'word_count']])
```

Listing 11.1: String manipulation in pandas

### 11.2 String Methods and Pattern Matching

```
1 # String splitting
2 text_data[['first_name', 'last_name']] = text_data['name'].str.split(' ', expand
3     =True)
4
5 # Pattern matching and extraction
6 text_data['has_python'] = text_data['description'].str.contains('python', case=
7     False)
8 text_data['has_experience'] = text_data['description'].str.contains(r'\d+\s+
9     years?', regex=True)
```

```

8 text_data['experience_years'] = text_data['description'].str.extract(r'(\d+)\s+
  years?')[0].fillna(0).astype(int)
9
10 # String replacement and cleaning
11 text_data['phone_clean'] = text_data['phone'].str.replace(r'\D', '', regex=True)
  # Remove non-digits
12 text_data['description_clean'] = text_data['description'].str.replace(r'[^\w\s]'
  , '', regex=True) # Remove punctuation
13
14 # String slicing and positioning
15 text_data['name_initials'] = text_data['first_name'].str[0] + text_data['
  last_name'].str[0]
16 text_data['domain_extension'] = text_data['email_domain'].str.split('.').str[-1]
17
18 print("After advanced string operations:")
19 print(text_data[['name', 'first_name', 'last_name', 'email_domain', 'has_python'
  , 'experience_years', 'phone_clean']])

```

Listing 11.2: Advanced string operations

## 11.3 Text Analysis and NLP Basics

```

1 import re
2 from collections import Counter
3
4 # Text analysis functions
5 def clean_text(text):
6     """Clean and preprocess text"""
7     text = text.lower()
8     text = re.sub(r'[^\w\s]', '', text) # Remove punctuation
9     text = re.sub(r'\d+', '', text) # Remove numbers
10    return text.strip()
11
12 def get_word_frequencies(series):
13     """Get word frequencies from a text series"""
14     all_words = ' '.join(series.apply(clean_text)).split()
15     return Counter(all_words)
16
17 # Apply text cleaning
18 text_data['description_cleaned'] = text_data['description'].apply(clean_text)
19
20 # Word frequency analysis
21 word_freq = get_word_frequencies(text_data['description'])
22 print("Most common words in descriptions:")
23 for word, count in word_freq.most_common(10):
24     print(f"{word}: {count}")
25
26 # TF-IDF like analysis (simplified)
27 def calculate_tf(description, word):
28     """Calculate term frequency for a word in description"""
29     words = clean_text(description).split()
30     return words.count(word) / len(words) if words else 0
31
32 # Analyze specific terms
33 terms = ['python', 'data', 'developer', 'machine', 'experience']
34 for term in terms:
35     text_data[f'tf_{term}'] = text_data['description'].apply(lambda x:
        calculate_tf(x, term))

```

```

36
37 print("\nTerm frequencies:")
38 print(text_data[['name'] + [f'tf_{term}' for term in terms]])

```

Listing 11.3: Text analysis with pandas

## 11.4 Regular Expressions with Pandas

```

1 # Complex pattern matching with regex
2 text_data['job_level'] = text_data['description'].str.extract(r'(senior|junior|
   lead|principal)', flags=re.IGNORECASE)[0].fillna('mid')
3 text_data['technologies'] = text_data['description'].str.findall(r'(python|
   machine learning|AI|web|backend|frontend|agile)', flags=re.IGNORECASE)
4 text_data['tech_count'] = text_data['technologies'].str.len()
5
6 # Email validation and parsing
7 text_data['is_valid_email'] = text_data['email'].str.match(r'^[a-zA-Z0-9._%+-]+@
   [a-zA-Z0-9.-]+\.[a-zA-Z]{2,}$')
8 text_data['email_username'] = text_data['email'].str.extract(r'^([^\@]+)\@')[0]
9
10 # Phone number formatting
11 def format_phone(phone):
12     """Standardize phone number format"""
13     digits = re.sub(r'\D', '', phone)
14     if len(digits) == 10:
15         return f"({digits[:3]}) {digits[3:6]}-{digits[6:]}"
16     return phone
17
18 text_data['phone_formatted'] = text_data['phone'].apply(format_phone)
19
20 # Extract skills using regex patterns
21 skill_patterns = {
22     'programming': r'\b(python|java|javascript|c\+\+|c#|ruby|go)\b',
23     'data_science': r'\b(machine learning|data science|AI|artificial
   intelligence|deep learning)\b',
24     'web': r'\b(web|frontend|backend|fullstack|react|angular|vue)\b',
25     'methodology': r'\b(agile|scrum|kanban|waterfall)\b'
26 }
27
28 for skill_category, pattern in skill_patterns.items():
29     text_data[f'has_{skill_category}'] = text_data['description'].str.contains(
   pattern, case=False, regex=True)
30
31 print("Regex extraction results:")
32 print(text_data[['name', 'job_level', 'tech_count', 'is_valid_email', '
   phone_formatted'] +
33         [f'has_{cat}' for cat in skill_patterns.keys()]])

```

Listing 11.4: Regular expression operations

## 11.5 Categorical Text Analysis

```

1 # Convert text to categorical when appropriate
2 text_data['job_level'] = text_data['job_level'].astype('category')
3 text_data['domain_extension'] = text_data['domain_extension'].astype('category')
4

```

```

5 # Analyze categorical text data
6 print("Job level distribution:")
7 print(text_data['job_level'].value_counts())
8
9 print("\nDomain extension distribution:")
10 print(text_data['domain_extension'].value_counts())
11
12 # Text similarity (basic)
13 from sklearn.feature_extraction.text import TfidfVectorizer
14 from sklearn.metrics.pairwise import cosine_similarity
15
16 # Calculate text similarity between descriptions
17 vectorizer = TfidfVectorizer(stop_words='english', max_features=50)
18 tfidf_matrix = vectorizer.fit_transform(text_data['description_cleaned'])
19 similarity_matrix = cosine_similarity(tfidf_matrix)
20
21 # Create similarity DataFrame
22 similarity_df = pd.DataFrame(
23     similarity_matrix,
24     index=text_data['name'],
25     columns=text_data['name']
26 )
27
28 print("\nText similarity matrix:")
29 print(similarity_df.round(3))
30
31 # Find most similar profiles
32 def find_similar_profiles(name, n=3):
33     """Find n most similar profiles to given name"""
34     similarities = similarity_df[name].sort_values(ascending=False)
35     return similarities.iloc[1:n+1] # Exclude self
36
37 print("\nSimilar profiles to 'John Doe':")
38 print(find_similar_profiles('John Doe'))

```

Listing 11.5: Categorical analysis of text data

## 11.6 Exercises

1. Clean and standardize a dataset of messy customer names and addresses.
2. Extract specific information from text using regular expressions (dates, prices, etc.).
3. Perform sentiment analysis on a set of product reviews.
4. Create word clouds and frequency analysis from text data.
5. Build a simple text classification system using pandas string operations.

## 11.7 Solutions

```

1.
1 # Clean customer data
2 customers = pd.DataFrame({
3     'name': ['john DOE ', 'jane SMITH ', ' bob j. ', 'ALICE-MARIE', '
4     'charlie wilson '],
5     'address': ['123 main st', '456 ELM AVENUE', '789 oak Rd.', '101 PINE
6     'BLVD', '202 MAPLE STREET'],

```

```

5     'phone': ['(123)456-7890', '987.654.3210', '555-123-4567', '444 555
      6666', '3332221111']
6 })
7
8 # Clean names
9 customers['name_clean'] = (customers['name']
10                             .str.strip()
11                             .str.title()
12                             .str.replace(r'\s+', ' ', regex=True))
13
14 # Clean addresses
15 customers['address_clean'] = (customers['address']
16                                 .str.strip()
17                                 .str.title()
18                                 .str.replace(r'\.', '', regex=True)
19                                 .str.replace(r'\s+', ' ', regex=True))
20
21 # Clean phones
22 customers['phone_clean'] = (customers['phone']
23                             .str.replace(r'\D', '', regex=True)
24                             .str.replace(r'(\d{3})(\d{3})(\d{4})', r'(\1
      \2-\3', regex=True))
25
26 print("Cleaned customer data:")
27 print(customers)

```

```

1 # Extract information from text
2 texts = pd.Series([
3     'Buy now for only $29.99! Limited offer until 2024-12-31',
4     'Price: $15.50, available from 2024-01-15 to 2024-06-30',
5     'Special discount: $99.00 valid until 2024-03-31',
6     'Cost: $45.75, promotion ends 2024-02-28'
7 ])
8
9 # Extract prices and dates
10 extracted_data = pd.DataFrame({
11     'original': texts,
12     'price': texts.str.extract(r'\$(\d+\.\d*)')[0].astype(float),
13     'start_date': texts.str.extract(r'(\d{4}-\d{2}-\d{2})')[0],
14     'end_date': texts.str.extract(r'(\d{4}-\d{2}-\d{2})')[1]
15 })
16
17 print("Extracted information:")
18 print(extracted_data)

```

3.

```

1 # Simple sentiment analysis
2 reviews = pd.DataFrame({
3     'review': [
4         'This product is amazing! I love it so much.',
5         'Terrible quality, very disappointed.',
6         'It is okay, nothing special.',
7         'Excellent product, would definitely recommend!',
8         'Poor customer service and bad product.'
9     ]
10 })
11
12 # Simple sentiment dictionary
13 positive_words = ['amazing', 'love', 'excellent', 'great', 'good', '
      recommend', 'awesome']

```

```

14 negative_words = ['terrible', 'disappointed', 'poor', 'bad', 'awful', '
    horrible']
15
16 def simple_sentiment(text):
17     text_lower = text.lower()
18     positive_count = sum(1 for word in positive_words if word in text_lower
19 )
19     negative_count = sum(1 for word in negative_words if word in text_lower
20 )
21
22     if positive_count > negative_count:
23         return 'positive'
24     elif negative_count > positive_count:
25         return 'negative'
26     else:
27         return 'neutral'
28
29 reviews['sentiment'] = reviews['review'].apply(simple_sentiment)
30 reviews['positive_words'] = reviews['review'].str.lower().str.findall('|'.
    join(positive_words))
31 reviews['negative_words'] = reviews['review'].str.lower().str.findall('|'.
    join(negative_words))
32
33 print("Sentiment analysis results:")
34 print(reviews)

```

```

1 # Word frequency analysis
2 from wordcloud import WordCloud
3 import matplotlib.pyplot as plt
4
5 # Sample text data
6 documents = [
7     'machine learning data science artificial intelligence',
8     'python programming data analysis pandas numpy',
9     'deep learning neural networks tensorflow pytorch',
10    'data visualization matplotlib seaborn plotly',
11    'statistics probability hypothesis testing'
12 ]
13
14 # Create word frequency
15 all_text = ' '.join(documents)
16 words = all_text.split()
17 word_freq = pd.Series(words).value_counts()
18
19 print("Word frequencies:")
20 print(word_freq)
21
22 # Create simple word cloud (text-based)
23 print("\nWord cloud (text):")
24 for word, count in word_freq.head(10).items():
25     print(f"{word}: {'#' * count}")
26
27 # Alternatively, create actual word cloud image
28 # wordcloud = WordCloud(width=800, height=400, background_color='white').
    generate(all_text)
29 # plt.figure(figsize=(10, 5))
30 # plt.imshow(wordcloud, interpolation='bilinear')
31 # plt.axis('off')
32 # plt.show()

```

```

4.
1 # Simple text classification
2 articles = pd.DataFrame({
3     'title': [
4         'Stock Market Reaches All Time High',
5         'Football Team Wins Championship Game',
6         'New Technology Company IPO Success',
7         'Basketball Player Sets New Record',
8         'Economic Growth Exceeds Expectations'
9     ],
10    'category': ['business', 'sports', 'business', 'sports', 'business']
11 })
12
13 # Create feature vectors based on keywords
14 business_keywords = ['stock', 'market', 'company', 'economic', 'growth', '
    ipo']
15 sports_keywords = ['football', 'basketball', 'championship', 'game', '
    player', 'record']
16
17 def extract_features(title):
18     title_lower = title.lower()
19     return {
20         'business_score': sum(1 for word in business_keywords if word in
            title_lower),
21         'sports_score': sum(1 for word in sports_keywords if word in
            title_lower),
22         'length': len(title.split())
23     }
24
25 # Apply feature extraction
26 features = articles['title'].apply(extract_features)
27 feature_df = pd.DataFrame(features.tolist())
28 articles_with_features = pd.concat([articles, feature_df], axis=1)
29
30 # Simple classification rule
31 def classify_article(row):
32     if row['business_score'] > row['sports_score']:
33         return 'business'
34     elif row['sports_score'] > row['business_score']:
35         return 'sports'
36     else:
37         return 'unknown'
38
39 articles_with_features['predicted_category'] = articles_with_features.apply
    (classify_article, axis=1)
40 articles_with_features['correct'] = articles_with_features['category'] ==
    articles_with_features['predicted_category']
41
42 print("Classification results:")
43 print(articles_with_features)
44 print(f"\nAccuracy: {articles_with_features['correct'].mean():.2f}")

```



# Chapter 12

## Advanced Visualization with Pandas

### 12.1 Basic Plotting with Pandas

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 # Create sample data for visualization
6 np.random.seed(42)
7 dates = pd.date_range('2024-01-01', periods=100, freq='D')
8 viz_data = pd.DataFrame({
9     'date': dates,
10    'sales': np.random.randint(1000, 5000, 100).cumsum(),
11    'temperature': 20 + 10 * np.sin(2 * np.pi * np.arange(100) / 30) + np.random
12        .normal(0, 3, 100),
13    'category': np.random.choice(['A', 'B', 'C'], 100),
14    'region': np.random.choice(['North', 'South', 'East', 'West'], 100)
15 })
16 viz_data = viz_data.set_index('date')
17 print("Visualization data:")
18 print(viz_data.head())
19
20 # Basic line plot
21 plt.figure(figsize=(12, 8))
22
23 plt.subplot(2, 2, 1)
24 viz_data['sales'].plot(title='Sales Over Time', color='blue')
25 plt.ylabel('Sales')
26
27 plt.subplot(2, 2, 2)
28 viz_data['temperature'].plot(title='Temperature Over Time', color='red')
29 plt.ylabel('Temperature')
30
31 plt.subplot(2, 2, 3)
32 viz_data['sales'].hist(bins=20, alpha=0.7, color='green')
33 plt.title('Sales Distribution')
34 plt.xlabel('Sales')
35
36 plt.subplot(2, 2, 4)
37 viz_data['category'].value_counts().plot(kind='bar', color='orange')
38 plt.title('Category Distribution')
39 plt.xlabel('Category')
40
41 plt.tight_layout()
42 plt.show()
```

Listing 12.1: Basic pandas visualization

## 12.2 Advanced Plot Types

```

1 # Create more complex visualizations
2 plt.figure(figsize=(15, 10))
3
4 # 1. Area plot for cumulative data
5 plt.subplot(2, 3, 1)
6 category_sales = viz_data.groupby([viz_data.index.month, 'category'])['sales'].
    sum().unstack()
7 category_sales.plot(kind='area', alpha=0.7, title='Monthly Sales by Category')
8 plt.ylabel('Sales')
9
10 # 2. Scatter plot with coloring
11 plt.subplot(2, 3, 2)
12 for region in viz_data['region'].unique():
13     region_data = viz_data[viz_data['region'] == region]
14     plt.scatter(region_data['temperature'], region_data['sales'],
15                 alpha=0.6, label=region, s=50)
16 plt.xlabel('Temperature')
17 plt.ylabel('Sales')
18 plt.title('Sales vs Temperature by Region')
19 plt.legend()
20
21 # 3. Box plot
22 plt.subplot(2, 3, 3)
23 viz_data.boxplot(column='sales', by='category', ax=plt.gca())
24 plt.title('Sales Distribution by Category')
25 plt.suptitle('') # Remove automatic title
26
27 # 4. Heatmap (correlation)
28 plt.subplot(2, 3, 4)
29 numeric_data = viz_data.select_dtypes(include=[np.number])
30 correlation_matrix = numeric_data.corr()
31 plt.imshow(correlation_matrix, cmap='coolwarm', aspect='auto')
32 plt.colorbar()
33 plt.xticks(range(len(correlation_matrix.columns)), correlation_matrix.columns,
34             rotation=45)
35 plt.yticks(range(len(correlation_matrix.columns)), correlation_matrix.columns)
36 plt.title('Correlation Heatmap')
37
38 # 5. Pie chart
39 plt.subplot(2, 3, 5)
40 region_counts = viz_data['region'].value_counts()
41 plt.pie(region_counts.values, labels=region_counts.index, autopct='%1.1f%%')
42 plt.title('Region Distribution')
43
44 # 6. Density plot
45 plt.subplot(2, 3, 6)
46 viz_data['sales'].plot(kind='density', color='purple')
47 plt.title('Sales Density Distribution')
48 plt.xlabel('Sales')
49
50 plt.tight_layout()
51 plt.show()

```

Listing 12.2: Advanced plotting techniques

## 12.3 Time Series Visualization

```

1 # Time series specific plots
2 plt.figure(figsize=(15, 10))
3
4 # 1. Time series with rolling average
5 plt.subplot(2, 2, 1)
6 viz_data['sales'].plot(alpha=0.5, label='Daily Sales', color='blue')
7 viz_data['sales'].rolling(7).mean().plot(label='7-Day Moving Average', color='
    red', linewidth=2)
8 plt.title('Sales with Moving Average')
9 plt.legend()
10
11 # 2. Seasonal decomposition visualization
12 from statsmodels.tsa.seasonal import seasonal_decompose
13
14 # Create a sample time series with seasonality
15 ts_with_seasonality = 100 + 5 * np.arange(100) + 10 * np.sin(2 * np.pi * np.
    arange(100) / 30)
16 decomposition = seasonal_decompose(ts_with_seasonality, model='additive', period
    =30)
17
18 plt.subplot(2, 2, 2)
19 decomposition.observed.plot(title='Original Series', color='blue')
20 plt.subplot(2, 2, 3)
21 decomposition.trend.plot(title='Trend', color='green')
22 plt.subplot(2, 2, 4)
23 decomposition.seasonal.plot(title='Seasonal', color='red')
24
25 plt.tight_layout()
26 plt.show()
27
28 # 3. Seasonal plot
29 plt.figure(figsize=(12, 6))
30 # Create monthly data for seasonal analysis
31 monthly_data = viz_data.resample('M').agg({'sales': 'sum', 'temperature': 'mean'
    })
32 monthly_data['month'] = monthly_data.index.month
33
34 # Pivot for seasonal plot
35 seasonal_pivot = monthly_data.pivot_table(values='sales', columns=monthly_data.
    index.year,
36                                           index='month', aggfunc='sum')
37 seasonal_pivot.plot(marker='o', title='Seasonal Sales Pattern by Year')
38 plt.xlabel('Month')
39 plt.ylabel('Sales')
40 plt.legend(title='Year')
41 plt.show()

```

Listing 12.3: Time series specific visualizations

## 12.4 Interactive Visualizations

```

1 try:
2     import plotly.express as px
3     import plotly.graph_objects as go
4     from plotly.subplots import make_subplots

```

```

5
6 # Create interactive line plot
7 fig = px.line(viz_data.reset_index(), x='date', y='sales',
8               title='Interactive Sales Over Time',
9               labels={'sales': 'Sales Amount', 'date': 'Date'})
10 fig.show()
11
12 # Create interactive scatter plot
13 fig2 = px.scatter(viz_data.reset_index(), x='temperature', y='sales',
14                  color='region', size='sales',
15                  hover_data=['category'],
16                  title='Sales vs Temperature by Region')
17 fig2.show()
18
19 # Create interactive bar chart
20 monthly_summary = viz_data.resample('M').agg({
21     'sales': 'sum',
22     'temperature': 'mean'
23 }).reset_index()
24
25 fig3 = go.Figure()
26 fig3.add_trace(go.Bar(x=monthly_summary['date'], y=monthly_summary['sales'],
27                       name='Sales', marker_color='blue'))
28 fig3.add_trace(go.Scatter(x=monthly_summary['date'], y=monthly_summary['
29     temperature'],
30                           name='Temperature', yaxis='y2', line=dict(color='
31     red'))))
32
33 fig3.update_layout(
34     title='Monthly Sales and Temperature',
35     xaxis=dict(title='Date'),
36     yaxis=dict(title='Sales', side='left'),
37     yaxis2=dict(title='Temperature', side='right', overlaying='y'),
38     legend=dict(x=0, y=1)
39 )
40 fig3.show()
41
42 except ImportError:
43     print("Plotly not installed. Install with: pip install plotly")
44
45 # Fallback to matplotlib
46 fig, ax1 = plt.subplots(figsize=(12, 6))
47
48 ax1.bar(monthly_summary['date'], monthly_summary['sales'], alpha=0.7, color=
49     'blue', label='Sales')
50 ax1.set_xlabel('Date')
51 ax1.set_ylabel('Sales', color='blue')
52 ax1.tick_params(axis='y', labelcolor='blue')
53
54 ax2 = ax1.twinx()
55 ax2.plot(monthly_summary['date'], monthly_summary['temperature'], color='red
56     ', marker='o', label='Temperature')
57 ax2.set_ylabel('Temperature', color='red')
58 ax2.tick_params(axis='y', labelcolor='red')
59
60 plt.title('Monthly Sales and Temperature')
61 plt.show()

```

Listing 12.4: Interactive plotting with Plotly

## 12.5 Custom Styling and Themes

```

1 # Apply custom styles
2 plt.style.use('seaborn-v0_8') # Use seaborn style
3
4 # Create customized plots
5 fig, ((ax1, ax2), (ax3, ax4)) = plt.subplots(2, 2, figsize=(15, 10))
6
7 # Customized line plot
8 ax1.plot(viz_data.index, viz_data['sales'],
9         color='#2E86AB', linewidth=2, marker='o', markersize=4)
10 ax1.set_title('Sales Trend', fontsize=14, fontweight='bold')
11 ax1.set_ylabel('Sales', fontsize=12)
12 ax1.grid(True, alpha=0.3)
13 ax1.tick_params(axis='x', rotation=45)
14
15 # Customized bar plot
16 category_performance = viz_data.groupby('category')['sales'].mean().sort_values
17     ()
18 colors = ['#A23B72', '#F18F01', '#C73E1D']
19 category_performance.plot(kind='barh', ax=ax2, color=colors)
20 ax2.set_title('Average Sales by Category', fontsize=14, fontweight='bold')
21 ax2.set_xlabel('Average Sales', fontsize=12)
22
23 # Customized scatter plot
24 scatter = ax3.scatter(viz_data['temperature'], viz_data['sales'],
25                     c=pd.factorize(viz_data['region'])[0],
26                     cmap='viridis', alpha=0.7, s=50)
27 ax3.set_title('Sales vs Temperature', fontsize=14, fontweight='bold')
28 ax3.set_xlabel('Temperature', fontsize=12)
29 ax3.set_ylabel('Sales', fontsize=12)
30 plt.colorbar(scatter, ax=ax3, label='Region')
31
32 # Customized box plot
33 box_data = [viz_data[viz_data['region'] == region]['sales'] for region in
34             viz_data['region'].unique()]
35 box_plot = ax4.boxplot(box_data, labels=viz_data['region'].unique(),
36                       patch_artist=True)
37
38 # Color the boxes
39 colors = ['#FF6B6B', '#4ECDC4', '#45B7D1', '#96CEB4']
40 for patch, color in zip(box_plot['boxes'], colors):
41     patch.set_facecolor(color)
42 ax4.set_title('Sales Distribution by Region', fontsize=14, fontweight='bold')
43 ax4.set_ylabel('Sales', fontsize=12)
44
45 plt.tight_layout()
46 plt.show()
47
48 # Save the figure
49 plt.savefig('custom_visualizations.png', dpi=300, bbox_inches='tight')

```

Listing 12.5: Customizing plot styles

## 12.6 Exercises

1. Create a comprehensive dashboard with multiple subplots showing different aspects of your data.

2. Build an interactive visualization that allows exploring relationships between variables.
3. Create a time series visualization with multiple moving averages and trend lines.
4. Design a custom color scheme and apply it to all your visualizations.
5. Compare different plot types for the same data and analyze which conveys information most effectively.

## 12.7 Solutions

```

1.
1 # Comprehensive dashboard
2 def create_dashboard(df):
3     fig, axes = plt.subplots(2, 3, figsize=(18, 12))
4
5     # 1. Sales trend
6     df['sales'].plot(ax=axes[0,0], color='blue', title='Sales Trend')
7     axes[0,0].set_ylabel('Sales')
8
9     # 2. Sales by category
10    df.groupby('category')['sales'].sum().plot(kind='bar', ax=axes[0,1],
11                                                color=['red', 'green', 'blue'],
12                                                title='Total Sales by
13                                                    Category')
14
15    # 3. Temperature distribution
16    df['temperature'].hist(ax=axes[0,2], bins=20, alpha=0.7, color='orange')
17    axes[0,2].set_title('Temperature Distribution')
18
19    # 4. Sales vs Temperature scatter
20    for category in df['category'].unique():
21        category_data = df[df['category'] == category]
22        axes[1,0].scatter(category_data['temperature'], category_data['
23            sales'],
24                           label=category, alpha=0.6)
25    axes[1,0].set_xlabel('Temperature')
26    axes[1,0].set_ylabel('Sales')
27    axes[1,0].set_title('Sales vs Temperature')
28    axes[1,0].legend()
29
30    # 5. Monthly sales
31    monthly_sales = df['sales'].resample('M').sum()
32    monthly_sales.plot(ax=axes[1,1], marker='o', color='purple')
33    axes[1,1].set_title('Monthly Sales')
34    axes[1,1].tick_params(axis='x', rotation=45)
35
36    # 6. Region distribution
37    df['region'].value_counts().plot(kind='pie', ax=axes[1,2], autopct='
38        %1.1f%%')
39    axes[1,2].set_title('Region Distribution')
40
41    plt.suptitle('Sales Data Dashboard', fontsize=16, fontweight='bold')
42    plt.tight_layout()
43    return fig
44
45 dashboard = create_dashboard(viz_data)

```

```
43 plt.show()
```

```
1 # Interactive exploration (simplified)
2 try:
3     import plotly.express as px
4
5     # Create interactive scatter plot
6     interactive_df = viz_data.reset_index()
7     fig = px.scatter(interactive_df, x='temperature', y='sales',
8                     color='region', size='sales',
9                     hover_data=['date', 'category'],
10                    title='Interactive Sales Exploration',
11                    labels={'sales': 'Sales', 'temperature': 'Temperature'})
12
13    # Add trend line
14    fig.add_traces(px.scatter(interactive_df, x='temperature', y='sales',
15                             trendline="lowess").data[1])
16
17    fig.show()
18
19 except ImportError:
20     # Fallback static version
21     plt.figure(figsize=(10, 6))
22     scatter = plt.scatter(viz_data['temperature'], viz_data['sales'],
23                         c=pd.factorize(viz_data['region'])[0],
24                         s=viz_data['sales']/100, alpha=0.6)
25     plt.colorbar(scatter, label='Region')
26     plt.xlabel('Temperature')
27     plt.ylabel('Sales')
28     plt.title('Sales vs Temperature by Region')
29     plt.show()
```

3.

```
1 # Time series with multiple moving averages
2 plt.figure(figsize=(12, 8))
3
4 # Original series
5 plt.subplot(2, 1, 1)
6 viz_data['sales'].plot(color='black', alpha=0.3, label='Daily Sales')
7
8 # Multiple moving averages
9 windows = [7, 14, 30]
10 colors = ['red', 'blue', 'green']
11 for window, color in zip(windows, colors):
12     viz_data['sales'].rolling(window).mean().plot(
13         label=f'{window}-Day MA', color=color, linewidth=2)
14
15 plt.title('Sales with Multiple Moving Averages')
16 plt.legend()
17 plt.ylabel('Sales')
18
19 # Residuals from trend
20 plt.subplot(2, 1, 2)
21 trend = viz_data['sales'].rolling(30).mean()
22 residuals = viz_data['sales'] - trend
23 residuals.plot(color='purple', alpha=0.7)
24 plt.axhline(y=0, color='red', linestyle='--', alpha=0.5)
25 plt.title('Residuals from 30-Day Moving Average')
```

```

26 plt.ylabel('Residuals')
27
28 plt.tight_layout()
29 plt.show()

```

```

1  # Custom color scheme
2  custom_colors = {
3      'primary': '#2E86AB',
4      'secondary': '#A23B72',
5      'accent1': '#F18F01',
6      'accent2': '#C73E1D',
7      'accent3': '#3C91E6'
8  }
9
10 # Apply custom colors
11 plt.figure(figsize=(12, 8))
12
13 # Plot 1: Line plot with custom colors
14 plt.subplot(2, 2, 1)
15 viz_data['sales'].plot(color=custom_colors['primary'], linewidth=2)
16 plt.title('Sales Trend', color=custom_colors['primary'])
17 plt.ylabel('Sales')
18
19 # Plot 2: Bar plot
20 plt.subplot(2, 2, 2)
21 category_data = viz_data.groupby('category')['sales'].mean()
22 category_data.plot(kind='bar', color=[custom_colors['secondary'],
23                                     custom_colors['accent1'],
24                                     custom_colors['accent2']])
25 plt.title('Sales by Category')
26 plt.xticks(rotation=45)
27
28 # Plot 3: Scatter plot
29 plt.subplot(2, 2, 3)
30 plt.scatter(viz_data['temperature'], viz_data['sales'],
31             alpha=0.6, color=custom_colors['accent3'])
32 plt.xlabel('Temperature')
33 plt.ylabel('Sales')
34 plt.title('Sales vs Temperature')
35
36 # Plot 4: Distribution
37 plt.subplot(2, 2, 4)
38 viz_data['sales'].hist(bins=20, alpha=0.7, color=custom_colors['accent1'])
39 plt.title('Sales Distribution')
40 plt.xlabel('Sales')
41
42 plt.suptitle('Custom Styled Visualizations', fontsize=16,
43             color=custom_colors['primary'], fontweight='bold')
44 plt.tight_layout()
45 plt.show()

```

4.

```

1  # Plot type comparison
2  comparison_data = viz_data.groupby('region')['sales'].agg(['mean', 'std', 'count'])
3
4  fig, axes = plt.subplots(2, 2, figsize=(15, 10))
5
6  # Bar plot

```



```

7 comparison_data['mean'].plot(kind='bar', ax=axes[0,0], color='skyblue')
8 axes[0,0].set_title('Bar Plot - Mean Sales by Region')
9 axes[0,0].tick_params(axis='x', rotation=45)
10
11 # Pie chart
12 axes[0,1].pie(comparison_data['mean'], labels=comparison_data.index,
13               autopct='%1.1f%%')
14 axes[0,1].set_title('Pie Chart - Mean Sales Distribution')
15
16 # Box plot
17 region_sales = [viz_data[viz_data['region'] == region]['sales'] for region
18                  in comparison_data.index]
19 axes[1,0].boxplot(region_sales, labels=comparison_data.index)
20 axes[1,0].set_title('Box Plot - Sales Distribution by Region')
21 axes[1,0].tick_params(axis='x', rotation=45)
22
23 # Violin plot (using boxplot as substitute if violin not available)
24 try:
25     axes[1,1].violinplot(region_sales, showmeans=True)
26     axes[1,1].set_xticks(range(1, len(comparison_data.index) + 1))
27     axes[1,1].set_xticklabels(comparison_data.index)
28     axes[1,1].set_title('Violin Plot - Sales Distribution')
29 except:
30     # Fallback to enhanced box plot
31     bp = axes[1,1].boxplot(region_sales, labels=comparison_data.index,
32                             patch_artist=True)
33     for patch in bp['boxes']:
34         patch.set_facecolor('lightgreen')
35     axes[1,1].set_title('Enhanced Box Plot - Sales Distribution')
36
37 plt.suptitle('Comparison of Plot Types for Regional Sales Data', fontsize
38             =16)
39 plt.tight_layout()
40 plt.show()
41
42 # Analysis
43 print("Effectiveness Analysis:")
44 print("Bar plot: Good for comparing exact values")
45 print("Pie chart: Good for showing proportions")
46 print("Box plot: Good for showing distribution statistics")
47 print("Violin plot: Best for detailed distribution visualization")

```

## Real-World Case Studies and Best Practices

### 13.1 Complete Data Analysis Workflow

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 def complete_analysis_project():
7     """Complete data analysis workflow"""
8     print("=== COMPLETE DATA ANALYSIS WORKFLOW ===\n")
9
10    # 1. Data Loading and Exploration
11    print("1. DATA LOADING AND EXPLORATION")
12    # Simulate real-world dataset
13    np.random.seed(42)
14    n_customers = 10000
15
16    customer_data = pd.DataFrame({
17        'customer_id': range(1, n_customers + 1),
18        'age': np.random.randint(18, 70, n_customers),
19        'income': np.random.normal(50000, 20000, n_customers).clip(20000,
20        150000),
21        'city': np.random.choice(['New York', 'Los Angeles', 'Chicago', 'Houston',
22        'Phoenix'],
23        n_customers, p=[0.2, 0.3, 0.15, 0.2, 0.15]),
24        'signup_date': pd.date_range('2020-01-01', periods=n_customers, freq='H')[:n_customers],
25        'loyalty_tier': np.random.choice(['Bronze', 'Silver', 'Gold', 'Platinum'],
26        n_customers, p=[0.4, 0.3, 0.2, 0.1])
27    })
28
29    # Add some missing values to simulate real data
30    customer_data.loc[customer_data.sample(frac=0.05).index, 'income'] = np.nan
31    customer_data.loc[customer_data.sample(frac=0.03).index, 'city'] = None
32
33    print("Dataset shape:", customer_data.shape)
34    print("\nFirst 5 rows:")
35    print(customer_data.head())
36    print("\nData types:")
37    print(customer_data.dtypes)
38    print("\nMissing values:")
39    print(customer_data.isnull().sum())
40
41    # 2. Data Cleaning and Preprocessing
42    print("\n2. DATA CLEANING AND PREPROCESSING")
43
44    # Handle missing values
```

```

43 customer_data_clean = customer_data.copy()
44 customer_data_clean['income'] = customer_data_clean['income'].fillna(
45     customer_data_clean.groupby('loyalty_tier')['income'].transform('median'
46 )
47 customer_data_clean['city'] = customer_data_clean['city'].fillna('Unknown')
48
49 # Feature engineering
50 customer_data_clean['signup_year'] = customer_data_clean['signup_date'].dt.
    year
51 customer_data_clean['signup_month'] = customer_data_clean['signup_date'].dt.
    month
52 customer_data_clean['age_group'] = pd.cut(customer_data_clean['age'],
53     bins=[18, 25, 35, 45, 55, 65, 70],
54     labels=['18-25', '26-35', '36-45', '
        46-55', '56-65', '66+'])
55 customer_data_clean['income_category'] = pd.qcut(customer_data_clean['income
    '],
56     q=4, labels=['Low', 'Medium',
        'High', 'Very High'])
57
58 print("After cleaning and feature engineering:")
59 print(customer_data_clean.info())
60
61 # 3. Exploratory Data Analysis
62 print("\n3. EXPLORATORY DATA ANALYSIS")
63
64 # Basic statistics
65 print("\nNumerical columns statistics:")
66 print(customer_data_clean[['age', 'income']].describe())
67
68 print("\nCategorical columns value counts:")
69 for col in ['city', 'loyalty_tier', 'age_group', 'income_category']:
70     print(f"\n{col}:")
71     print(customer_data_clean[col].value_counts())
72
73 # 4. Advanced Analysis
74 print("\n4. ADVANCED ANALYSIS")
75
76 # Customer segmentation analysis
77 segmentation = customer_data_clean.groupby(['loyalty_tier', 'age_group']).
    agg({
78     'customer_id': 'count',
79     'income': ['mean', 'median', 'std'],
80     'age': 'mean'
81 }).round(2)
82
83 print("\nCustomer Segmentation Analysis:")
84 print(segmentation)
85
86 # Time-based analysis
87 monthly_signups = customer_data_clean.groupby(
88     [customer_data_clean['signup_date'].dt.to_period('M'), 'loyalty_tier']
89 )['customer_id'].count().unstack().fillna(0)
90
91 print("\nMonthly Signups by Loyalty Tier:")
92 print(monthly_signups.tail())
93
94 # 5. Visualization
95 print("\n5. DATA VISUALIZATION")

```

```

96 plt.figure(figsize=(15, 10))
97
98
99 # Subplot 1: Age distribution
100 plt.subplot(2, 3, 1)
101 customer_data_clean['age'].hist(bins=20, alpha=0.7, color='skyblue')
102 plt.title('Age Distribution')
103 plt.xlabel('Age')
104 plt.ylabel('Count')
105
106 # Subplot 2: Income by loyalty tier
107 plt.subplot(2, 3, 2)
108 customer_data_clean.boxplot(column='income', by='loyalty_tier', ax=plt.gca()
109 )
110 plt.title('Income Distribution by Loyalty Tier')
111 plt.suptitle('')
112
113 # Subplot 3: City distribution
114 plt.subplot(2, 3, 3)
115 customer_data_clean['city'].value_counts().plot(kind='bar', color='
116 lightgreen')
117 plt.title('Customer Distribution by City')
118 plt.xticks(rotation=45)
119
120 # Subplot 4: Monthly signups trend
121 plt.subplot(2, 3, 4)
122 monthly_signups.plot(ax=plt.gca(), marker='o')
123 plt.title('Monthly Signups Trend')
124 plt.xlabel('Month')
125 plt.ylabel('Number of Signups')
126 plt.legend(title='Loyalty Tier')
127
128 # Subplot 5: Age group by loyalty tier
129 plt.subplot(2, 3, 5)
130 cross_tab = pd.crosstab(customer_data_clean['age_group'],
131                          customer_data_clean['loyalty_tier'],
132                          normalize='index')
133 cross_tab.plot(kind='bar', ax=plt.gca())
134 plt.title('Loyalty Tier Distribution by Age Group')
135 plt.xlabel('Age Group')
136 plt.ylabel('Proportion')
137 plt.legend(title='Loyalty Tier')
138 plt.xticks(rotation=45)
139
140 # Subplot 6: Income vs Age scatter
141 plt.subplot(2, 3, 6)
142 for tier in customer_data_clean['loyalty_tier'].unique():
143     tier_data = customer_data_clean[customer_data_clean['loyalty_tier'] ==
144     tier]
145     plt.scatter(tier_data['age'], tier_data['income'], alpha=0.6, label=tier
146     , s=20)
147 plt.xlabel('Age')
148 plt.ylabel('Income')
149 plt.title('Income vs Age by Loyalty Tier')
150 plt.legend()
151
152 plt.tight_layout()
153 plt.show()
154
155 # 6. Insights and Recommendations

```

```

152 print("\n6. INSIGHTS AND RECOMMENDATIONS")
153
154 # Calculate key metrics
155 avg_income_by_tier = customer_data_clean.groupby('loyalty_tier')['income'].
    mean()
156 customer_count_by_city = customer_data_clean['city'].value_counts()
157 monthly_growth = monthly_signups.sum(axis=1).pct_change().mean()
158
159 print(f"\nKey Insights:")
160 print(f"- Average income increases with loyalty tier: {avg_income_by_tier.
    to_dict()}")
161 print(f"- Top city for customers: {customer_count_by_city.index[0]} ({
    customer_count_by_city.iloc[0]} customers)")
162 print(f"- Average monthly growth rate: {monthly_growth:.2%}")
163
164 # Business recommendations
165 print(f"\nBusiness Recommendations:")
166 print("1. Focus retention efforts on Platinum and Gold tier customers (
    highest value)")
167 print("2. Develop targeted marketing for high-income age groups (45-65)")
168 print("3. Expand presence in underpenetrated cities")
169 print("4. Create loyalty programs to upgrade Bronze/Silver customers")
170
171 return customer_data_clean, segmentation
172
173 # Run the complete analysis
174 final_data, segmentation_results = complete_analysis_project()

```

Listing 13.1: End-to-end data analysis project

## 13.2 Performance Optimization Case Study

```

1 def performance_optimization_demo():
2     """Demonstrate performance optimization techniques"""
3     print("=== PERFORMANCE OPTIMIZATION CASE STUDY ===\n")
4
5     # Create a large dataset for performance testing
6     print("Creating large dataset...")
7     n_rows = 1000000
8     n_cols = 20
9
10    large_data = pd.DataFrame(
11        np.random.randn(n_rows, n_cols),
12        columns=[f'feature_{i}' for i in range(n_cols)]
13    )
14
15    # Add categorical columns
16    large_data['category'] = np.random.choice(['A', 'B', 'C', 'D'], n_rows)
17    large_data['subcategory'] = np.random.choice(['X', 'Y', 'Z'], n_rows)
18    large_data['flag'] = np.random.choice([True, False], n_rows)
19
20    print(f"Dataset size: {n_rows:,} rows x {n_cols} columns")
21    print(f"Memory usage: {large_data.memory_usage(deep=True).sum() / 1024**2:.2
        f} MB")
22
23    # 1. Memory Optimization
24    print("\n1. MEMORY OPTIMIZATION")
25

```

```

26 def optimize_memory(df):
27     """Optimize DataFrame memory usage"""
28     optimized_df = df.copy()
29
30     # Downcast numeric columns
31     for col in optimized_df.select_dtypes(include=['integer']):
32         optimized_df[col] = pd.to_numeric(optimized_df[col], downcast='
            integer')
33
34     for col in optimized_df.select_dtypes(include=['float']):
35         optimized_df[col] = pd.to_numeric(optimized_df[col], downcast='float
            ')
36
37     # Convert object columns to category when appropriate
38     for col in optimized_df.select_dtypes(include=['object']):
39         num_unique = optimized_df[col].nunique()
40         num_total = len(optimized_df[col])
41         if num_unique / num_total < 0.5:
42             optimized_df[col] = optimized_df[col].astype('category')
43
44     return optimized_df
45
46 import time
47 start_time = time.time()
48 optimized_data = optimize_memory(large_data)
49 optimization_time = time.time() - start_time
50
51 original_memory = large_data.memory_usage(deep=True).sum()
52 optimized_memory = optimized_data.memory_usage(deep=True).sum()
53
54 print(f"Original memory: {original_memory / 1024**2:.2f} MB")
55 print(f"Optimized memory: {optimized_memory / 1024**2:.2f} MB")
56 print(f"Memory reduction: {(1 - optimized_memory/original_memory)*100:.1f}%"
57     )
58 print(f"Optimization time: {optimization_time:.2f} seconds")
59
60 # 2. Operation Performance Comparison
61 print("\n2. OPERATION PERFORMANCE COMPARISON")
62
63 # Test different grouping methods
64 operations = []
65
66 # Method 1: Basic groupby
67 start_time = time.time()
68 result1 = large_data.groupby('category').mean()
69 time1 = time.time() - start_time
70 operations.append(('Basic groupby', time1))
71
72 # Method 2: Optimized groupby
73 start_time = time.time()
74 result2 = optimized_data.groupby('category').mean()
75 time2 = time.time() - start_time
76 operations.append(('Optimized groupby', time2))
77
78 # Method 3: Using categorical
79 start_time = time.time()
80 large_data_cat = large_data.copy()
81 large_data_cat['category'] = large_data_cat['category'].astype('category')
82 result3 = large_data_cat.groupby('category').mean()
83 time3 = time.time() - start_time

```

```

83     operations.append(('Categorical groupby', time3))
84
85     # Print performance results
86     print("\nGroupBy Performance:")
87     for operation, duration in operations:
88         print(f"{operation}: {duration:.4f} seconds")
89
90     # 3. Vectorization vs Apply
91     print("\n3. VECTORIZATION VS APPLY")
92
93     # Complex calculation: weighted sum of features
94     def complex_calc_apply(row):
95         return sum(row[f'feature_{i}'] * (i + 1) for i in range(10))
96
97     def complex_calc_vectorized(df):
98         return sum(df[f'feature_{i}'] * (i + 1) for i in range(10))
99
100    # Test apply
101    start_time = time.time()
102    result_apply = large_data.iloc[:10000].apply(complex_calc_apply, axis=1)
103    time_apply = time.time() - start_time
104
105    # Test vectorized
106    start_time = time.time()
107    result_vectorized = complex_calc_vectorized(large_data.iloc[:10000])
108    time_vectorized = time.time() - start_time
109
110    print(f"Apply method: {time_apply:.4f} seconds")
111    print(f"Vectorized method: {time_vectorized:.4f} seconds")
112    print(f"Speedup: {time_apply/time_vectorized:.1f}x")
113
114    # 4. Chunk Processing for Very Large Data
115    print("\n4. CHUNK PROCESSING DEMO")
116
117    def process_in_chunks(df, chunk_size=10000):
118        """Process DataFrame in chunks"""
119        chunks = []
120        total_rows = len(df)
121
122        for i in range(0, total_rows, chunk_size):
123            chunk = df.iloc[i:i + chunk_size]
124            # Simulate complex processing
125            processed_chunk = chunk[chunk['flag'] == True].groupby('category').mean()
126            chunks.append(processed_chunk)
127
128            if (i // chunk_size) % 10 == 0:
129                print(f"Processed {min(i + chunk_size, total_rows):,} rows...")
130
131        # Combine results
132        final_result = pd.concat(chunks).groupby(level=0).mean()
133        return final_result
134
135    start_time = time.time()
136    chunk_result = process_in_chunks(large_data)
137    chunk_time = time.time() - start_time
138
139    print(f"Chunk processing time: {chunk_time:.2f} seconds")
140    print(f"Chunk processing result shape: {chunk_result.shape}")
141

```

```

142     return optimized_data, operations
143
144 # Run performance optimization demo
145 optimized_df, performance_stats = performance_optimization_demo()

```

Listing 13.2: Performance optimization in practice

## 13.3 Best Practices and Code Quality

```

1 def demonstrate_best_practices():
2     """Show pandas best practices and coding standards"""
3     print("=== PANDAS BEST PRACTICES ===\n")
4
5     # 1. Code Organization
6     print("1. CODE ORGANIZATION AND READABILITY")
7
8     # Good: Clean, readable code
9     def calculate_customer_metrics(customers_df, orders_df):
10         """
11         Calculate key customer metrics from customer and order data.
12
13         Parameters:
14         customers_df (DataFrame): Customer demographic data
15         orders_df (DataFrame): Order transaction data
16
17         Returns:
18         DataFrame: Customer metrics with lifetime value and frequency
19         """
20         # Merge customer and order data
21         customer_orders = pd.merge(
22             customers_df,
23             orders_df,
24             on='customer_id',
25             how='left'
26         )
27
28         # Calculate metrics using method chaining
29         customer_metrics = (
30             customer_orders
31             .groupby('customer_id')
32             .agg({
33                 'order_amount': ['sum', 'mean', 'count'],
34                 'order_date': ['min', 'max']
35             })
36             .round(2)
37         )
38
39         # Flatten column names
40         customer_metrics.columns = [
41             'total_spent', 'avg_order_value', 'order_count',
42             'first_order', 'last_order'
43         ]
44
45         # Calculate additional metrics
46         customer_metrics['lifetime_days'] = (
47             customer_metrics['last_order'] - customer_metrics['first_order']
48         ).dt.days
49

```



```

50     return customer_metrics
51
52 # 2. Error Handling and Validation
53 print("\n2. ERROR HANDLING AND VALIDATION")
54
55 def safe_data_processing(df, required_columns=None):
56     """
57     Safely process DataFrame with validation and error handling.
58     """
59     if required_columns is None:
60         required_columns = ['customer_id', 'order_amount']
61
62     # Input validation
63     if not isinstance(df, pd.DataFrame):
64         raise TypeError("Input must be a pandas DataFrame")
65
66     missing_columns = [col for col in required_columns if col not in df.
67                        columns]
68     if missing_columns:
69         raise ValueError(f"Missing required columns: {missing_columns}")
70
71     try:
72         # Create a copy to avoid modifying original data
73         processed_df = df.copy()
74
75         # Data validation and cleaning
76         processed_df = processed_df.dropna(subset=required_columns)
77
78         # Type validation
79         if 'order_amount' in processed_df.columns:
80             processed_df['order_amount'] = pd.to_numeric(
81                 processed_df['order_amount'], errors='coerce'
82             )
83             # Remove invalid values
84             processed_df = processed_df[processed_df['order_amount'] > 0]
85
86         return processed_df
87
88     except Exception as e:
89         print(f"Error during processing: {e}")
90         return pd.DataFrame() # Return empty DataFrame on error
91
92 # 3. Efficient Data Processing Patterns
93 print("\n3. EFFICIENT DATA PROCESSING PATTERNS")
94
95 def efficient_data_pipeline(raw_data):
96     """
97     Demonstrate efficient data processing pipeline.
98     """
99     return (
100         raw_data
101         # Initial filtering
102         .query('order_amount > 0 and customer_id.notna()')
103         # Type conversion
104         .assign(
105             order_date=lambda x: pd.to_datetime(x['order_date']),
106             order_amount=lambda x: pd.to_numeric(x['order_amount'])
107         )
108         # Feature engineering
109         .assign(

```

```

109         order_year=lambda x: x['order_date'].dt.year,
110         order_month=lambda x: x['order_date'].dt.month,
111         amount_category=lambda x: pd.cut(
112             x['order_amount'],
113             bins=[0, 50, 200, 500, float('inf')],
114             labels=['Small', 'Medium', 'Large', 'Very Large']
115         )
116     )
117     # Grouping and aggregation
118     .groupby(['customer_id', 'order_year'])
119     .agg({
120         'order_amount': ['sum', 'mean', 'count'],
121         'order_date': ['min', 'max']
122     })
123     # Flatten columns
124     .pipe(lambda df: df.set_axis(
125         ['total_amount', 'avg_amount', 'order_count',
126          'first_order', 'last_order'],
127         axis=1
128     ))
129     # Final calculations
130     .assign(
131         lifetime_days=lambda x: (x['last_order'] - x['first_order']).dt.
132             days,
133         daily_spend=lambda x: x['total_amount'] / x['lifetime_days'].
134             clip(1)
135     )
136     .round(2)
137
138 # 4. Memory and Performance Considerations
139 print("\n4. MEMORY AND PERFORMANCE CONSIDERATIONS")
140
141 def memory_efficient_processing(df):
142     """
143     Process data with memory efficiency in mind.
144     """
145     # Sample data for memory optimization
146     result = {}
147
148     # Use appropriate dtypes from start
149     dtypes = {
150         'customer_id': 'int32',
151         'order_amount': 'float32',
152         'category': 'category'
153     }
154
155     # Select only needed columns
156     needed_columns = ['customer_id', 'order_amount', 'order_date', 'category']
157     efficient_df = df[needed_columns].copy()
158
159     # Convert to efficient dtypes
160     for col, dtype in dtypes.items():
161         if col in efficient_df.columns:
162             efficient_df[col] = efficient_df[col].astype(dtype)
163
164     # Process in chunks if very large
165     chunk_results = []
166     chunk_size = 10000

```

```

166     for i in range(0, len(efficient_df), chunk_size):
167         chunk = efficient_df.iloc[i:i + chunk_size]
168
169         # Process chunk
170         chunk_result = (
171             chunk
172             .groupby('category')
173             .agg({'order_amount': 'sum'})
174         )
175         chunk_results.append(chunk_result)
176
177     # Combine results
178     final_result = pd.concat(chunk_results).groupby(level=0).sum()
179
180     result['memory_usage'] = efficient_df.memory_usage(deep=True).sum()
181     result['processing_time'] = "Simulated"
182     result['final_result'] = final_result
183
184     return result
185
186 # 5. Testing and Validation
187 print("\n5. TESTING AND VALIDATION")
188
189 def create_test_data():
190     """Create test data for validation"""
191     return pd.DataFrame({
192         'customer_id': [1, 2, 3, 4, 5],
193         'order_amount': [100.0, 200.0, 150.0, 300.0, 250.0],
194         'order_date': pd.date_range('2024-01-01', periods=5),
195         'category': ['A', 'B', 'A', 'C', 'B']
196     })
197
198 def test_data_processing():
199     """Test the data processing function"""
200     test_df = create_test_data()
201
202     try:
203         result = efficient_data_pipeline(test_df)
204
205         # Validation checks
206         assert len(result) > 0, "Result should not be empty"
207         assert 'total_amount' in result.columns, "Should have total_amount column"
208         assert result['total_amount'].sum() > 0, "Total amount should be positive"
209
210         print(" All tests passed!")
211         return True
212
213     except Exception as e:
214         print(f" Test failed: {e}")
215         return False
216
217 # Run the test
218 test_passed = test_data_processing()
219
220 print(f"\nSummary:")
221 print("- Write clean, documented code")
222 print("- Use method chaining for readability")
223

```

```

224     print("- Implement proper error handling")
225     print("- Optimize for memory and performance")
226     print("- Always test your code")
227     print("- Use appropriate data types")
228     print("- Process large datasets in chunks")
229
230     return test_passed
231
232 # Demonstrate best practices
233 best_practices_result = demonstrate_best_practices()

```

Listing 13.3: Pandas best practices

## 13.4 Final Project: Complete Business Analysis

```

1 def final_business_analysis_project():
2     """
3     Complete business analysis project demonstrating all pandas skills.
4     """
5     print("=== FINAL BUSINESS ANALYSIS PROJECT ===\n")
6
7     # 1. Problem Statement
8     print("1. PROBLEM STATEMENT")
9     print("Analyze e-commerce data to identify:")
10    print("- Customer segmentation and behavior patterns")
11    print("- Sales trends and seasonality")
12    print("- Product performance and inventory insights")
13    print("- Revenue optimization opportunities")
14
15    # 2. Data Generation (Simulating real business data)
16    print("\n2. DATA GENERATION AND LOADING")
17
18    np.random.seed(42)
19    n_transactions = 50000
20    n_customers = 5000
21    n_products = 200
22
23    # Generate customer data
24    customers = pd.DataFrame({
25        'customer_id': range(1, n_customers + 1),
26        'join_date': pd.date_range('2020-01-01', periods=n_customers, freq='H')
27        [:n_customers],
28        'region': np.random.choice(['North', 'South', 'East', 'West', 'International'],
29                                   n_customers, p=[0.25, 0.25, 0.2, 0.2, 0.1]),
30        'customer_tier': np.random.choice(['New', 'Regular', 'VIP'],
31                                           n_customers, p=[0.6, 0.35, 0.05]),
32        'age': np.random.randint(18, 70, n_customers),
33        'lifetime_value': np.random.exponential(1000, n_customers).clip(10,
34                                10000)
35    })
36
37    # Generate product catalog
38    products = pd.DataFrame({
39        'product_id': range(1, n_products + 1),
40        'product_name': [f'Product_{i:03d}' for i in range(1, n_products + 1)],
41        'category': np.random.choice(['Electronics', 'Clothing', 'Home', 'Books',
42                                     'Sports'],

```

```

40         n_products, p=[0.3, 0.25, 0.2, 0.15, 0.1]),
41         'subcategory': np.random.choice(['A', 'B', 'C', 'D'], n_products),
42         'cost_price': np.random.uniform(5, 500, n_products),
43         'retail_price': lambda x: x['cost_price'] * np.random.uniform(1.2, 3.0,
44             n_products)
45     })
46     products['retail_price'] = products['cost_price'] * np.random.uniform(1.2,
47         3.0, n_products)
48     products['margin'] = products['retail_price'] - products['cost_price']
49
50     # Generate transaction data
51     transactions = pd.DataFrame({
52         'transaction_id': range(1, n_transactions + 1),
53         'customer_id': np.random.randint(1, n_customers + 1, n_transactions),
54         'product_id': np.random.randint(1, n_products + 1, n_transactions),
55         'transaction_date': pd.date_range('2023-01-01', periods=n_transactions,
56             freq='H')[:n_transactions],
57         'quantity': np.random.randint(1, 5, n_transactions),
58         'discount_pct': np.random.choice([0, 0.1, 0.15, 0.2, 0.3],
59             n_transactions,
60             p=[0.5, 0.2, 0.15, 0.1, 0.05])
61     })
62
63     # Merge data and calculate final amounts
64     transactions = transactions.merge(products[['product_id', 'retail_price']],
65         on='product_id')
66     transactions['final_price'] = transactions['retail_price'] * (1 -
67         transactions['discount_pct'])
68     transactions['revenue'] = transactions['final_price'] * transactions['
69         quantity']
70
71     print(f"Customers: {len(customers):,}")
72     print(f"Products: {len(products):,}")
73     print(f"Transactions: {len(transactions):,}")
74
75     # 3. Comprehensive Analysis
76     print("\n3. COMPREHENSIVE ANALYSIS")
77
78     # 3.1 Customer Analysis
79     print("\n3.1 CUSTOMER ANALYSIS")
80
81     customer_behavior = transactions.merge(customers, on='customer_id')
82
83     customer_metrics = (
84         customer_behavior
85         .groupby('customer_id')
86         .agg({
87             'revenue': ['sum', 'mean', 'count'],
88             'transaction_date': ['min', 'max', 'nunique'],
89             'quantity': 'sum',
90             'customer_tier': 'first',
91             'region': 'first'
92         })
93     )
94     .pipe(lambda df: df.set_axis([
95         'total_revenue', 'avg_transaction', 'transaction_count',
96         'first_purchase', 'last_purchase', 'active_days',
97         'total_quantity', 'tier', 'region'
98     ], axis=1))
99     .assign(

```

```

92         lifetime_days=lambda x: (x['last_purchase'] - x['first_purchase']).
           dt.days,
93         daily_revenue=lambda x: x['total_revenue'] / x['lifetime_days'].clip
           (1),
94         customer_type=lambda x: np.where(
95             x['total_revenue'] > x['total_revenue'].quantile(0.8),
96             'High Value',
97             np.where(x['total_revenue'] > x['total_revenue'].quantile(0.5),
98                     'Medium Value', 'Low Value')
99         )
100     )
101
102     print("Customer Segments:")
103     print(customer_metrics['customer_type'].value_counts())
104     print(f"\nAverage Revenue by Tier:")
105     print(customer_metrics.groupby('tier')['total_revenue'].mean().round(2))
106
107     # 3.2 Product Analysis
108     print("\n3.2 PRODUCT ANALYSIS")
109
110     product_performance = (
111         transactions
112         .merge(products, on='product_id')
113         .groupby(['product_id', 'product_name', 'category'])
114         .agg({
115             'revenue': 'sum',
116             'quantity': 'sum',
117             'transaction_id': 'count',
118             'retail_price': 'mean',
119             'discount_pct': 'mean'
120         })
121         .rename(columns={'transaction_id': 'units_sold'})
122         .assign(
123             avg_selling_price=lambda x: x['revenue'] / x['quantity'],
124             profit_margin=lambda x: (x['avg_selling_price'] -
125                                     products.set_index('product_id').loc[x.index.
126                                         get_level_values('product_id'), '
127                                         cost_price'].values) /
128                                     x['avg_selling_price']
129         )
130         .sort_values('revenue', ascending=False)
131     )
132
133     print("Top 10 Products by Revenue:")
134     print(product_performance.head(10)[['revenue', 'units_sold', 'profit_margin'
135                                         ]].round(2))
136
137     # 3.3 Time Series Analysis
138     print("\n3.3 TIME SERIES ANALYSIS")
139
140     # Daily revenue trends
141     daily_revenue = (
142         transactions
143         .set_index('transaction_date')
144         .resample('D')
145         .agg({'revenue': 'sum', 'transaction_id': 'count'})
146         .rename(columns={'transaction_id': 'daily_transactions'})
147     )

```

```

146 # Add rolling averages
147 daily_revenue['revenue_7d_ma'] = daily_revenue['revenue'].rolling(7).mean()
148 daily_revenue['revenue_30d_ma'] = daily_revenue['revenue'].rolling(30).mean()
149
150 # Seasonal analysis
151 daily_revenue['day_of_week'] = daily_revenue.index.dayofweek
152 daily_revenue['month'] = daily_revenue.index.month
153 daily_revenue['is_weekend'] = daily_revenue['day_of_week'].isin([5, 6])
154
155 print("Weekly Pattern (Average Revenue by Day):")
156 weekly_pattern = daily_revenue.groupby('day_of_week')['revenue'].mean()
157 print(weekly_pattern.round(2))
158
159 # 4. Advanced Insights
160 print("\n4. ADVANCED INSIGHTS AND RECOMMENDATIONS")
161
162 # Customer Lifetime Value prediction (simplified)
163 clv_data = customer_metrics.copy()
164 clv_data['predicted_6m_value'] = clv_data['daily_revenue'] * 180
165
166 # Product-category analysis
167 category_analysis = (
168     product_performance
169     .groupby('category')
170     .agg({
171         'revenue': 'sum',
172         'units_sold': 'sum',
173         'profit_margin': 'mean'
174     })
175     .sort_values('revenue', ascending=False)
176 )
177
178 # Regional performance
179 regional_performance = (
180     customer_metrics
181     .groupby('region')
182     .agg({
183         'total_revenue': 'sum',
184         'customer_id': 'count',
185         'daily_revenue': 'mean'
186     })
187     .rename(columns={'customer_id': 'customer_count'})
188     .assign(
189         revenue_per_customer=lambda x: x['total_revenue'] / x['customer_count']
190     )
191 )
192
193 print("\nKey Business Insights:")
194 print("=" * 50)
195
196 # Insight 1: Top performing categories
197 top_category = category_analysis.index[0]
198 top_category_revenue = category_analysis.iloc[0]['revenue']
199 print(f"1. {top_category} is the highest revenue category (${top_category_revenue:,.0f})")
200
201 # Insight 2: Best customer segment

```

```

202     best_segment = customer_metrics.groupby('customer_type')['total_revenue'].
        mean().idxmax()
203     best_segment_value = customer_metrics.groupby('customer_type')['
        total_revenue'].mean().max()
204     print(f"2. {best_segment} customers generate ${best_segment_value:,.0f} on
        average")
205
206     # Insight 3: Regional performance
207     best_region = regional_performance['revenue_per_customer'].idxmax()
208     best_region_value = regional_performance['revenue_per_customer'].max()
209     print(f"3. {best_region} has highest revenue per customer (${
        best_region_value:,.0f})")
210
211     # Insight 4: Time-based insights
212     best_day = weekly_pattern.idxmax()
213     days = ['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday',
        'Sunday']
214     print(f"4. {days[best_day]} is the best performing day of the week")
215
216     print("\nStrategic Recommendations:")
217     print("=" * 50)
218     print("1. Focus marketing budget on high-performing categories and regions")
219     print("2. Develop retention programs for high-value customer segments")
220     print("3. Optimize inventory for weekend and peak day demand")
221     print("4. Create targeted promotions for underperforming regions")
222     print("5. Implement customer loyalty programs to increase lifetime value")
223
224     # 5. Export Results
225     print("\n5. RESULTS EXPORT")
226
227     # Create summary report
228     summary_report = {
229         'total_revenue': transactions['revenue'].sum(),
230         'total_customers': customers['customer_id'].nunique(),
231         'total_products': products['product_id'].nunique(),
232         'avg_transaction_value': transactions['revenue'].mean(),
233         'top_category': top_category,
234         'best_region': best_region,
235         'customer_acquisition_cost': 50, # Example metric
236         'customer_lifetime_value': clv_data['predicted_6m_value'].mean()
237     }
238
239     summary_df = pd.DataFrame([summary_report]).T
240     summary_df.columns = ['Value']
241
242     print("\nBusiness Summary:")
243     print(summary_df)
244
245     return {
246         'customers': customers,
247         'products': products,
248         'transactions': transactions,
249         'customer_metrics': customer_metrics,
250         'product_performance': product_performance,
251         'daily_revenue': daily_revenue,
252         'summary': summary_df
253     }
254
255     # Run the final project
256     final_results = final_business_analysis_project()

```



```

257
258 print("\n" + "="*60)
259 print("CONGRATULATIONS! You've completed the Pandas Mastery course!")
260 print("="*60)
261 print("\nYou now have comprehensive skills in:")
262 print("  Data manipulation and cleaning")
263 print("  Advanced analysis and aggregation")
264 print("  Time series analysis")
265 print("  Performance optimization")
266 print("  Real-world business applications")
267 print("  Best practices and professional coding")
268 print("\nContinue practicing with real datasets and explore:")
269 print("- Machine learning integration")
270 print("- Big data technologies (Dask, PySpark)")
271 print("- Database integration (SQL, NoSQL)")
272 print("- Web scraping and API data collection")
273 print("- Advanced visualization libraries")

```

Listing 13.4: Final comprehensive project

## 13.5 Exercises

1. Implement the complete business analysis project with your own simulated data.
2. Identify three additional business metrics that could be valuable and calculate them.
3. Create a comprehensive dashboard visualizing all key insights from the analysis.
4. Optimize the code for better performance and memory usage.
5. Extend the analysis to include predictive modeling for customer behavior.

## 13.6 Solutions

```

1.
1  # Custom business analysis implementation
2  def custom_business_analysis():
3      # Generate custom business data
4      np.random.seed(123)
5
6      # Your custom data generation here
7      custom_customers = pd.DataFrame({
8          'customer_id': range(1000),
9          'segment': np.random.choice(['A', 'B', 'C'], 1000),
10         # Add more fields as needed
11     })
12
13     # Implement your analysis pipeline
14     results = {}
15
16     # Add your analysis steps
17     results['customer_summary'] = custom_customers.describe()
18
19     return results
20
21 # Run custom analysis
22 custom_results = custom_business_analysis()

```

```

1 # Additional business metrics
2 def calculate_additional_metrics(transactions, customers):
3     additional_metrics = {}
4
5     # 1. Customer retention rate
6     monthly_customers = transactions.groupby(
7         transactions['transaction_date'].dt.to_period('M')
8     )['customer_id'].nunique()
9     retention_rate = monthly_customers.pct_change().mean()
10    additional_metrics['retention_rate'] = retention_rate
11
12    # 2. Average order value by segment
13    segment_aov = transactions.merge(customers, on='customer_id').groupby('
14        customer_tier')['revenue'].mean()
15    additional_metrics['segment_aov'] = segment_aov
16
17    # 3. Product return rate (simulated)
18    # In real scenario, you'd have return data
19    additional_metrics['estimated_return_rate'] = 0.05 # 5% assumption
20
21    return additional_metrics
22
23    additional_metrics = calculate_additional_metrics(
24        final_results['transactions'],
25        final_results['customers']
26    )
27    print("Additional Metrics:", additional_metrics)

```

3.

```

1 # Comprehensive dashboard
2 def create_business_dashboard(results):
3     fig, axes = plt.subplots(3, 2, figsize=(15, 12))
4
5     # 1. Revenue trend
6     results['daily_revenue']['revenue'].plot(ax=axes[0,0], title='Daily
7         Revenue')
8
9     # 2. Customer segments
10    results['customer_metrics']['customer_type'].value_counts().plot(
11        kind='pie', ax=axes[0,1], title='Customer Segments'
12    )
13
14    # 3. Category performance
15    results['product_performance'].groupby('category')['revenue'].sum().
16    plot(
17        kind='bar', ax=axes[1,0], title='Revenue by Category'
18    )
19
20    # 4. Regional performance
21    regional_data = results['customer_metrics'].groupby('region')['
22        total_revenue'].sum()
23    regional_data.plot(kind='bar', ax=axes[1,1], title='Revenue by Region')
24
25    # 5. Price distribution
26    results['transactions']['final_price'].hist(ax=axes[2,0], bins=30,
27        alpha=0.7)
28    axes[2,0].set_title('Price Distribution')
29
30    # 6. Customer lifetime value
31    results['customer_metrics']['total_revenue'].plot(

```

```

28         kind='box', ax=axes[2,1], title='Customer Lifetime Value
           Distribution'
29     )
30
31     plt.suptitle('Business Performance Dashboard', fontsize=16, fontweight=
           'bold')
32     plt.tight_layout()
33     plt.show()
34
35 create_business_dashboard(final_results)

```

```

1  # Performance optimization
2  def optimize_business_analysis(transactions, customers, products):
3      # Use efficient dtypes
4      transactions_opt = transactions.astype({
5          'customer_id': 'int32',
6          'product_id': 'int32',
7          'quantity': 'int16',
8          'discount_pct': 'float32',
9          'revenue': 'float32'
10     })
11
12     customers_opt = customers.astype({
13         'customer_id': 'int32',
14         'region': 'category',
15         'customer_tier': 'category',
16         'age': 'int16',
17         'lifetime_value': 'float32'
18     })
19
20     # Process in chunks for large datasets
21     def process_chunked_analysis(df, chunk_size=10000):
22         chunks = []
23         for i in range(0, len(df), chunk_size):
24             chunk = df.iloc[i:i+chunk_size]
25             # Perform analysis on chunk
26             chunk_analysis = chunk.groupby('customer_id')['revenue'].sum()
27             chunks.append(chunk_analysis)
28         return pd.concat(chunks).groupby(level=0).sum()
29
30     # Use vectorized operations
31     customer_metrics_opt = (
32         transactions_opt
33         .merge(customers_opt, on='customer_id')
34         .groupby('customer_id')
35         .agg({
36             'revenue': 'sum',
37             'transaction_date': ['min', 'max'],
38             'quantity': 'sum'
39         })
40     )
41
42     return {
43         'memory_saved': f"{{(transactions.memory_usage().sum() -
           transactions_opt.memory_usage().sum()) / 1024**2:.1f}} MB",
44         'customer_metrics': customer_metrics_opt
45     }
46
47 optimization_results = optimize_business_analysis(
48     final_results['transactions'],

```

```

49     final_results['customers'],
50     final_results['products']
51 )
52 print("Optimization Results:", optimization_results['memory_saved'])

```

4.

```

1  # Predictive modeling extension
2  def customer_behavior_prediction(customer_metrics):
3      from sklearn.ensemble import RandomForestRegressor
4      from sklearn.model_selection import train_test_split
5      from sklearn.metrics import mean_absolute_error, r2_score
6
7      # Prepare features for prediction
8      features = customer_metrics.select_dtypes(include=[np.number]).copy()
9      features = features.fillna(features.mean())
10
11     # Create target variable (future revenue - simplified)
12     # In reality, you'd need temporal data for proper prediction
13     features['target'] = features['total_revenue'] * np.random.uniform(0.8,
14         1.2, len(features))
15
16     X = features.drop('target', axis=1)
17     y = features['target']
18
19     # Split data
20     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=
21         0.2, random_state=42)
22
23     # Train model
24     model = RandomForestRegressor(n_estimators=100, random_state=42)
25     model.fit(X_train, y_train)
26
27     # Predictions
28     y_pred = model.predict(X_test)
29
30     # Evaluation
31     mae = mean_absolute_error(y_test, y_pred)
32     r2 = r2_score(y_test, y_pred)
33
34     feature_importance = pd.DataFrame({
35         'feature': X.columns,
36         'importance': model.feature_importances_
37     }).sort_values('importance', ascending=False)
38
39     return {
40         'model': model,
41         'mae': mae,
42         'r2_score': r2,
43         'feature_importance': feature_importance,
44         'predictions': y_pred
45     }
46
47     # Note: This is a simplified example. Real predictive modeling would
48     # require
49     # proper temporal validation and feature engineering
50     try:
51         prediction_results = customer_behavior_prediction(final_results['
52             customer_metrics'])
53         print("Prediction Model Performance:")
54         print(f"MAE: {prediction_results['mae']:.2f}")

```

```
51     print(f"R2 Score: {prediction_results['r2_score']:.2f}")
52     print("\nTop Features:")
53     print(prediction_results['feature_importance'].head())
54 except Exception as e:
55     print(f"Prediction modeling skipped: {e}")
```

# Chapter 14

## Working with Different Data Sources

### 14.1 Database Operations with Pandas

```
1 import pandas as pd
2 import sqlite3
3 from sqlalchemy import create_engine
4
5 # SQLite database operations
6 def sqlite_operations():
7     # Create a connection
8     conn = sqlite3.connect('example.db')
9
10    # Read data from SQL table
11    df = pd.read_sql_query('SELECT * FROM customers', conn)
12
13    # Write DataFrame to SQL table
14    df.to_sql('new_table', conn, if_exists='replace', index=False)
15
16    # Close connection
17    conn.close()
18    return df
19
20 # Using SQLAlchemy for more database types
21 def sqlalchemy_operations():
22     # Create engine for different databases
23     # SQLite
24     engine_sqlite = create_engine('sqlite:///database.db')
25
26     # PostgreSQL
27     # engine_pg = create_engine('postgresql://user:password@localhost:5432/
28     #     database')
29
30     # MySQL
31     # engine_mysql = create_engine('mysql://user:password@localhost:3306/
32     #     database')
33
34     # Read from database
35     df = pd.read_sql_table('table_name', engine_sqlite)
36
37     # Read with custom query
38     query = """
39     SELECT customer_id, SUM(amount) as total_spent
40     FROM transactions
41     WHERE date >= '2024-01-01'
42     GROUP BY customer_id
43     """
44     df_custom = pd.read_sql_query(query, engine_sqlite)
45
46     # Write to database
47     df.to_sql('output_table', engine_sqlite, if_exists='replace', index=False)
```

```

47     return df
48
49 # Advanced database operations
50 def advanced_db_operations():
51     engine = create_engine('sqlite:///example.db')
52
53     # Read in chunks for large tables
54     chunk_size = 1000
55     chunks = []
56     for chunk in pd.read_sql_query('SELECT * FROM large_table',
57                                   engine, chunksize=chunk_size):
58         processed_chunk = chunk[chunk['value'] > 0]
59         chunks.append(processed_chunk)
60
61     large_df = pd.concat(chunks, ignore_index=True)
62
63     # Parameterized queries
64     params = {'min_amount': 100, 'start_date': '2024-01-01'}
65     safe_query = """
66     SELECT * FROM transactions
67     WHERE amount > %(min_amount)s
68     AND date >= %(start_date)s
69     """
70     df_safe = pd.read_sql_query(safe_query, engine, params=params)
71
72     return large_df, df_safe

```

Listing 14.1: Working with databases

## 14.2 Web APIs and JSON Data

```

1 import requests
2 import json
3
4 def api_operations():
5     # Basic API request
6     response = requests.get('https://api.example.com/data')
7
8     if response.status_code == 200:
9         # Convert JSON response to DataFrame
10         data = response.json()
11         df = pd.DataFrame(data['results'])
12
13         # Handle nested JSON
14         normalized_df = pd.json_normalize(data,
15                                           record_path='users',
16                                           meta=['page_info'])
17
18         return df, normalized_df
19
20     return pd.DataFrame()
21
22 # Working with complex JSON structures
23 def complex_json_operations():
24     # Sample complex JSON
25     complex_json = {
26         "company": "Tech Corp",
27         "employees": [

```

[illegible]



```

88         print("Users DataFrame shape:", users_df.shape)
89         print("Addresses DataFrame shape:", addresses_df.shape)
90
91         return users_df, addresses_df, companies_df
92
93     except Exception as e:
94         print(f"API error: {e}")
95         return pd.DataFrame(), pd.DataFrame(), pd.DataFrame()
96

```

Listing 14.2: Working with APIs and JSON

## 14.3 Web Scraping with Pandas

```

1  import requests
2  from bs4 import BeautifulSoup
3
4  def web_scraping_basics():
5      # Scrape HTML tables from web pages
6      url = 'https://example.com/table-page'
7
8      try:
9          response = requests.get(url)
10         response.raise_for_status()
11
12         # Parse HTML
13         soup = BeautifulSoup(response.content, 'html.parser')
14
15         # Find all tables
16         tables = soup.find_all('table')
17
18         # Read tables into pandas
19         dataframes = []
20         for i, table in enumerate(tables):
21             df = pd.read_html(str(table))[0]
22             df['table_source'] = f'table_{i+1}'
23             dataframes.append(df)
24
25         # Combine all tables
26         combined_df = pd.concat(dataframes, ignore_index=True)
27         return combined_df
28
29     except Exception as e:
30         print(f"Web scraping error: {e}")
31         return pd.DataFrame()
32
33 # Reading HTML tables directly
34 def read_html_tables():
35     # Read tables from HTML file or URL
36     try:
37         # From URL
38         dfs = pd.read_html('https://en.wikipedia.org/wiki/
39                             List_of_countries_by_GDP')
40
41         # From HTML string or file
42         html_string = """
43         <table>
44             <tr><th>Name</th><th>Age</th><th>City</th></tr>

```

```

44         <tr><td>John</td><td>25</td><td>New York</td></tr>
45         <tr><td>Jane</td><td>30</td><td>London</td></tr>
46     </table>
47     """
48
49     dfs = pd.read_html(html_string)
50     return dfs[0] if dfs else pd.DataFrame()
51
52 except Exception as e:
53     print(f"HTML reading error: {e}")
54     return pd.DataFrame()
55
56 # Advanced web scraping with data cleaning
57 def advanced_web_scraping():
58     # Simulated scraping with data cleaning
59     scraped_data = []
60
61     # Simulate multiple pages
62     for page in range(1, 4):
63         # In real scenario, this would be actual web requests
64         page_data = {
65             'product_name': f'Product {page}',
66             'price': 100 * page,
67             'rating': 4.0 + (page * 0.1),
68             'reviews': page * 10,
69             'page_number': page
70         }
71         scraped_data.append(page_data)
72
73     df = pd.DataFrame(scraped_data)
74
75     # Data cleaning operations
76     df['price_category'] = pd.cut(df['price'],
77                                 bins=[0, 150, 300, float('inf')],
78                                 labels=['Budget', 'Mid-range', 'Premium'])
79
80     df['popularity_score'] = df['rating'] * df['reviews']
81
82     return df

```

Listing 14.3: Web scraping and HTML parsing

## 14.4 Working with Excel and Office Documents

```

1 def advanced_excel_operations():
2     # Reading multiple sheets
3     excel_file = 'data.xlsx'
4
5     # Read all sheets
6     all_sheets = pd.read_excel(excel_file, sheet_name=None)
7
8     # Read specific sheets
9     specific_sheets = pd.read_excel(excel_file, sheet_name=['Sheet1', 'Sheet2'])
10
11     # Read with specific parameters
12     df_custom = pd.read_excel(excel_file,
13                               sheet_name=0,
14                               header=0,

```

```

15         usecols='A:D,F',
16         skiprows=[0, 2],
17         nrows=1000,
18         dtype={'column1': str, 'column2': float})
19
20 # Writing to Excel with formatting options
21 with pd.ExcelWriter('output.xlsx', engine='xlsxwriter') as writer:
22     df_custom.to_excel(writer, sheet_name='Processed_Data', index=False)
23
24     # Get workbook and worksheet objects
25     workbook = writer.book
26     worksheet = writer.sheets['Processed_Data']
27
28     # Add formatting (requires xlsxwriter)
29     header_format = workbook.add_format({
30         'bold': True,
31         'text_wrap': True,
32         'valign': 'top',
33         'fg_color': '#D7E4BC',
34         'border': 1
35     })
36
37     # Write headers with format
38     for col_num, value in enumerate(df_custom.columns.values):
39         worksheet.write(0, col_num, value, header_format)
40
41     return all_sheets, df_custom
42
43 # Working with multiple Excel files
44 def multiple_excel_files():
45     import glob
46
47     # Find all Excel files in directory
48     excel_files = glob.glob('*.xlsx')
49
50     all_data = []
51     for file in excel_files:
52         # Read each file
53         df = pd.read_excel(file)
54         df['source_file'] = file
55         all_data.append(df)
56
57     # Combine all data
58     combined_data = pd.concat(all_data, ignore_index=True)
59
60     # Summary by source file
61     summary = combined_data.groupby('source_file').agg({
62         'source_file': 'count',
63         **{col: 'mean' for col in combined_data.select_dtypes(include=[np.number])
64            }.columns}
65     ).rename(columns={'source_file': 'row_count'})
66
67     return combined_data, summary

```

Listing 14.4: Advanced Excel operations

## 14.5 Cloud Storage and Big Data

```

1 def cloud_storage_operations():
2     # Note: These are examples - actual implementation requires cloud SDKs
3
4     # Amazon S3
5     def s3_operations():
6         # import boto3
7         # s3 = boto3.client('s3')
8
9         # Read from S3
10        # df = pd.read_csv('s3://bucket-name/path/to/file.csv')
11
12        # Write to S3
13        # df.to_csv('s3://bucket-name/path/to/output.csv', index=False)
14
15        print("S3 operations would require boto3 installation")
16        return pd.DataFrame()
17
18    # Google Cloud Storage
19    def gcs_operations():
20        # from google.cloud import storage
21        # client = storage.Client()
22
23        # Read from GCS
24        # df = pd.read_csv('gs://bucket-name/path/to/file.csv')
25
26        print("GCS operations would require google-cloud-storage installation")
27        return pd.DataFrame()
28
29    # Azure Blob Storage
30    def azure_operations():
31        # from azure.storage.blob import BlobServiceClient
32        # Connect and read/write operations
33
34        print("Azure operations would require azure-storage-blob installation")
35        return pd.DataFrame()
36
37    return s3_operations(), gcs_operations(), azure_operations()
38
39 # Working with large datasets using Dask
40 def large_dataset_operations():
41     try:
42         import dask.dataframe as dd
43
44         # Create large dataset for demonstration
45         large_data = pd.DataFrame({
46             'id': range(1000000),
47             'value1': np.random.randn(1000000),
48             'value2': np.random.randn(1000000),
49             'category': np.random.choice(['A', 'B', 'C', 'D'], 1000000)
50         })
51
52         # Save to CSV for Dask demonstration
53         large_data.to_csv('large_dataset.csv', index=False)
54
55         # Read with Dask
56         ddf = dd.read_csv('large_dataset.csv')
57
58         # Perform operations (lazy evaluation)
59         result = (ddf[ddf['value1'] > 0]
60                 .groupby('category'))

```

```

61         .agg({'value1': 'mean', 'value2': 'std'})
62         .compute() # Actually compute the result
63
64     print("Dask processing completed")
65     return result
66
67 except ImportError:
68     print("Dask not installed. Install with: pip install dask")
69
70     # Fallback to pandas with chunk processing
71     chunks = []
72     chunk_size = 10000
73
74     for chunk in pd.read_csv('large_dataset.csv', chunksize=chunk_size):
75         processed_chunk = chunk[chunk['value1'] > 0]
76         chunks.append(processed_chunk)
77
78     result = pd.concat(chunks, ignore_index=True)
79     result = result.groupby('category').agg({
80         'value1': 'mean',
81         'value2': 'std'
82     })
83
84     return result

```

Listing 14.5: Working with cloud storage

## 14.6 Real-time Data and Streaming

```

1 def streaming_data_simulation():
2     # Simulate streaming data
3     import time
4     from datetime import datetime
5
6     streaming_data = []
7
8     for i in range(100):
9         # Simulate real-time data point
10         data_point = {
11             'timestamp': datetime.now(),
12             'value': np.random.randn(),
13             'sensor_id': f'sensor_{i % 5}',
14             'batch_number': i // 20
15         }
16         streaming_data.append(data_point)
17
18         # Simulate time delay
19         time.sleep(0.01)
20
21     # Convert to DataFrame
22     stream_df = pd.DataFrame(streaming_data)
23
24     # Real-time analytics
25     latest_values = stream_df.groupby('sensor_id').last()
26     rolling_avg = (stream_df.set_index('timestamp')
27                    .groupby('sensor_id')['value']
28                    .rolling('5s')
29                    .mean())

```

```

30         .reset_index())
31
32     return stream_df, latest_values, rolling_avg
33
34 # Data pipeline simulation
35 def data_pipeline_simulation():
36     # Simulate complete data pipeline
37     print("1. Data Extraction")
38     raw_data = pd.DataFrame({
39         'user_id': range(1000),
40         'action': np.random.choice(['click', 'view', 'purchase'], 1000),
41         'timestamp': pd.date_range('2024-01-01', periods=1000, freq='T'),
42         'value': np.random.randint(1, 100, 1000)
43     })
44
45     print("2. Data Transformation")
46     transformed_data = (raw_data
47         .assign(
48             hour=lambda x: x['timestamp'].dt.hour,
49             date=lambda x: x['timestamp'].dt.date,
50             value_category=lambda x: pd.cut(x['value'],
51                 bins=[0, 25, 50, 75, 100],
52                 labels=['Low', 'Medium', 'High', 'Very High'])
53         )
54     .query('value > 10')
55     )
56
57     print("3. Data Loading")
58     # Simulate loading to database or file
59     transformed_data.to_csv('processed_data.csv', index=False)
60
61     print("4. Data Analysis")
62     analysis = transformed_data.groupby(['date', 'action']).agg({
63         'user_id': 'count',
64         'value': ['mean', 'sum']
65     }).round(2)
66
67     return raw_data, transformed_data, analysis

```

Listing 14.6: Handling real-time data

## 14.7 Exercises

1. Connect to a SQLite database, create a table, and perform basic CRUD operations using pandas.
2. Fetch data from a public API and transform the JSON response into a structured DataFrame.
3. Scrape data from a webpage containing tables and clean the extracted data.
4. Create an Excel file with multiple sheets and different formatting options.
5. Build a complete data pipeline that extracts, transforms, and loads data from multiple sources.

## 14.8 Solutions

```

1.
1 def sqlite_exercise():
2     import sqlite3
3
4     # Create database and table
5     conn = sqlite3.connect('exercise.db')
6
7     # Sample data
8     employees = pd.DataFrame({
9         'id': range(1, 6),
10        'name': ['Alice', 'Bob', 'Charlie', 'Diana', 'Eva'],
11        'department': ['HR', 'Engineering', 'Engineering', 'Marketing', 'HR'],
12        'salary': [50000, 80000, 75000, 60000, 55000],
13        'join_date': pd.date_range('2020-01-01', periods=5, freq='M')
14    })
15
16    # Write to database
17    employees.to_sql('employees', conn, if_exists='replace', index=False)
18
19    # Read with query
20    query = "SELECT * FROM employees WHERE salary > 60000"
21    high_earners = pd.read_sql_query(query, conn)
22
23    # Update data
24    cursor = conn.cursor()
25    cursor.execute("UPDATE employees SET salary = 85000 WHERE name = 'Bob'")
26    conn.commit()
27
28    # Read updated data
29    updated_employees = pd.read_sql_query("SELECT * FROM employees", conn)
30
31    conn.close()
32    return employees, high_earners, updated_employees
33
34 emp_df, high_earners_df, updated_df = sqlite_exercise()
35 print("High earners:", high_earners_df)

```

```

1 def api_exercise():
2     import requests
3
4     try:
5         # Use JSONPlaceholder API
6         response = requests.get('https://jsonplaceholder.typicode.com/posts')
7
8         if response.status_code == 200:
9             posts_data = response.json()
10
11            # Convert to DataFrame
12            posts_df = pd.DataFrame(posts_data)
13
14            # Get user data
15            users_response = requests.get('https://jsonplaceholder.typicode.com/users')
16            users_data = users_response.json()
17            users_df = pd.json_normalize(users_data)

```

```

18
19     # Merge posts with user information
20     merged_df = posts_df.merge(users_df[['id', 'name', 'email']],
21                                left_on='userId', right_on='id',
22                                suffixes=('_post', '_user'))
23
24     # Analysis
25     posts_per_user = merged_df.groupby('name').agg({
26         'id_post': 'count',
27         'title': lambda x: ', '.join(x[:2]) # First two titles
28     }).rename(columns={'id_post': 'post_count', 'title': '
29         sample_titles'})
30
31     return merged_df, posts_per_user
32
33 except Exception as e:
34     print(f"API exercise error: {e}")
35     return pd.DataFrame(), pd.DataFrame()
36
37 posts_df, user_summary = api_exercise()
38 print("Posts per user:", user_summary)

```

3.

```

1 def web_scraping_exercise():
2     # For this exercise, we'll create HTML content and scrape it
3     html_content = """
4     <html>
5         <body>
6             <table border="1">
7                 <tr>
8                     <th>Product</th>
9                     <th>Price</th>
10                    <th>Rating</th>
11                    <th>Stock</th>
12                </tr>
13                <tr>
14                    <td>Laptop</td>
15                    <td>$999.99</td>
16                    <td>4.5</td>
17                    <td>15</td>
18                </tr>
19                <tr>
20                    <td>Mouse</td>
21                    <td>$29.99</td>
22                    <td>4.2</td>
23                    <td>50</td>
24                </tr>
25                <tr>
26                    <td>Keyboard</td>
27                    <td>$79.99</td>
28                    <td>4.7</td>
29                    <td>25</td>
30                </tr>
31            </table>
32        </body>
33    </html>
34    """
35
36    # Read HTML table
37    dfs = pd.read_html(html_content)

```



```

38     products_df = dfs[0] if dfs else pd.DataFrame()
39
40     # Data cleaning
41     if not products_df.empty:
42         products_df['Price'] = products_df['Price'].str.replace('$', '').
43             astype(float)
44         products_df['Rating'] = products_df['Rating'].astype(float)
45         products_df['Stock'] = products_df['Stock'].astype(int)
46
47         # Add calculated columns
48         products_df['Value_Score'] = (products_df['Rating'] * products_df['
49             Price'] / 100).round(2)
50         products_df['Stock_Status'] = np.where(products_df['Stock'] > 30, '
51             High',
52             np.where(products_df['Stock']
53                 > 10, 'Medium', 'Low'))
54
55     return products_df
56
57 cleaned_products = web_scraping_exercise()
58 print("Cleaned products data:")
59 print(cleaned_products)

```

```

1 def excel_exercise():
2     # Create sample data for different departments
3     sales_data = pd.DataFrame({
4         'Month': ['Jan', 'Feb', 'Mar', 'Apr', 'May'],
5         'Revenue': [10000, 12000, 11000, 13000, 14000],
6         'Expenses': [8000, 8500, 8200, 9000, 9500],
7         'Profit': [2000, 3500, 2800, 4000, 4500]
8     })
9
10    hr_data = pd.DataFrame({
11        'Employee': ['Alice', 'Bob', 'Charlie', 'Diana'],
12        'Department': ['Sales', 'Engineering', 'Engineering', 'Marketing'],
13        'Salary': [50000, 80000, 75000, 60000],
14        'Join_Date': pd.to_datetime(['2020-01-15', '2019-03-20', '
15            2021-07-10', '2020-11-05'])
16    })
17
18    inventory_data = pd.DataFrame({
19        'Product': ['Laptop', 'Mouse', 'Keyboard', 'Monitor'],
20        'Category': ['Electronics', 'Accessories', 'Accessories', '
21            Electronics'],
22        'Stock': [50, 200, 150, 30],
23        'Price': [999.99, 29.99, 79.99, 299.99],
24        'Last_Restock': pd.to_datetime(['2024-01-10', '2024-01-15', '
25            2024-01-12', '2024-01-08'])
26    })
27
28    # Write to Excel with different sheets
29    with pd.ExcelWriter('company_data.xlsx', engine='xlsxwriter') as writer:
30        # Sales data sheet
31        sales_data.to_excel(writer, sheet_name='Sales', index=False)
32
33        # HR data sheet
34        hr_data.to_excel(writer, sheet_name='HR', index=False)
35
36        # Inventory data sheet

```

```

34     inventory_data.to_excel(writer, sheet_name='Inventory', index=False
35     )
36
37     # Get workbook
38     workbook = writer.book
39
40     # Add formats
41     money_format = workbook.add_format({'num_format': '$#,##0.00'})
42     date_format = workbook.add_format({'num_format': 'yyyy-mm-dd'})
43     header_format = workbook.add_format({
44         'bold': True,
45         'bg_color': '#366092',
46         'font_color': 'white'
47     })
48
49     # Apply formats to each sheet
50     for sheet_name in ['Sales', 'HR', 'Inventory']:
51         worksheet = writer.sheets[sheet_name]
52
53         # Format headers
54         for col_num, value in enumerate(writer.sheets[sheet_name].table
55             .values[0]):
56             worksheet.write(0, col_num, value, header_format)
57
58         # Auto-adjust column widths
59         worksheet.set_column('A:Z', 15)
60
61     print("Excel file created: company_data.xlsx")
62     return sales_data, hr_data, inventory_data
63
64 sales, hr, inventory = excel_exercise()

```

4.

```

1 def data_pipeline_exercise():
2     print("=== COMPLETE DATA PIPELINE EXERCISE ===\n")
3
4     # Step 1: Extract from multiple sources
5     print("1. EXTRACTING DATA...")
6
7     # Source 1: CSV file
8     try:
9         customer_data = pd.DataFrame({
10             'customer_id': [1, 2, 3, 4, 5],
11             'name': ['Alice', 'Bob', 'Charlie', 'Diana', 'Eva'],
12             'join_date': pd.date_range('2023-01-01', periods=5, freq='M'),
13             'region': ['North', 'South', 'North', 'East', 'West']
14         })
15     except:
16         customer_data = pd.DataFrame()
17
18     # Source 2: Simulated API data
19     try:
20         orders_data = pd.DataFrame({
21             'order_id': [101, 102, 103, 104, 105],
22             'customer_id': [1, 2, 1, 3, 4],
23             'order_date': pd.date_range('2024-01-01', periods=5, freq='D'),
24             'amount': [150.0, 200.0, 75.0, 300.0, 125.0],
25             'status': ['completed', 'completed', 'pending', 'completed', '
26             shipped']
27         })

```

```

27     except:
28         orders_data = pd.DataFrame()
29
30     # Source 3: Simulated database data
31     try:
32         products_data = pd.DataFrame({
33             'product_id': [201, 202, 203, 204],
34             'product_name': ['Laptop', 'Mouse', 'Keyboard', 'Monitor'],
35             'category': ['Electronics', 'Accessories', 'Accessories', 'Electronics'],
36             'price': [999.99, 29.99, 79.99, 299.99]
37         })
38     except:
39         products_data = pd.DataFrame()
40
41     print(f"Extracted: {len(customer_data)} customers, {len(orders_data)} orders, {len(products_data)} products")
42
43     # Step 2: Transform data
44     print("\n2. TRANSFORMING DATA...")
45
46     # Clean and merge data
47     if not orders_data.empty and not customer_data.empty:
48         # Merge orders with customer data
49         orders_enriched = orders_data.merge(customer_data, on='customer_id')
50
51         # Add calculated fields
52         orders_enriched['order_weekday'] = orders_enriched['order_date'].dt.day_name()
53         orders_enriched['amount_category'] = pd.cut(orders_enriched['amount'],
54                                                     bins=[0, 100, 200, float('inf')],
55                                                     labels=['Small', 'Medium', 'Large'])
56
57         # Aggregate data
58         customer_summary = orders_enriched.groupby(['customer_id', 'name', 'region']).agg({
59             'order_id': 'count',
60             'amount': ['sum', 'mean'],
61             'order_date': ['min', 'max']
62         }).round(2)
63
64         # Flatten column names
65         customer_summary.columns = ['total_orders', 'total_spent', 'avg_order_value',
66                                     'first_order', 'last_order']
67
68         customer_summary = customer_summary.reset_index()
69
70         # Regional analysis
71         regional_summary = orders_enriched.groupby('region').agg({
72             'order_id': 'count',
73             'amount': ['sum', 'mean'],
74             'customer_id': 'nunique'
75         }).round(2)
76
77         regional_summary.columns = ['total_orders', 'total_revenue', '

```

```
78         'avg_order_value', 'unique_customers']
79     regional_summary = regional_summary.reset_index()
80
81     # Step 3: Load data
82     print("\n3. LOADING DATA...")
83
84     # Save to different formats
85     if 'orders_enriched' in locals():
86         orders_enriched.to_csv('processed_orders.csv', index=False)
87         customer_summary.to_csv('customer_analysis.csv', index=False)
88         regional_summary.to_csv('regional_analysis.csv', index=False)
89
90     # Save to Excel with multiple sheets
91     with pd.ExcelWriter('business_intelligence.xlsx') as writer:
92         orders_enriched.to_excel(writer, sheet_name='Orders_Detailed',
93                                 index=False)
94         customer_summary.to_excel(writer, sheet_name='Customer_Summary',
95                                 index=False)
96         regional_summary.to_excel(writer, sheet_name='Regional_Analysis',
97                                 index=False)
98
99     print("Pipeline completed successfully!")
100
101     # Return results
102     results = {
103         'customer_data': customer_data,
104         'orders_data': orders_data,
105         'products_data': products_data
106     }
107
108     if 'orders_enriched' in locals():
109         results.update({
110             'orders_enriched': orders_enriched,
111             'customer_summary': customer_summary,
112             'regional_summary': regional_summary
113         })
114
115     return results
116
117 # Run the pipeline
118 pipeline_results = data_pipeline_exercise()
```

# Glossary

**pandas** A Python library for data manipulation and analysis providing data structures and operations for manipulating numerical tables and time series.

**DataFrame** Two-dimensional labeled data structure with columns of potentially different types, similar to a spreadsheet or SQL table.

**Series** One-dimensional labeled array capable of holding any data type (integers, strings, floating point numbers, Python objects, etc.).

**NaN** Not a Number, a special floating-point value used to represent missing or undefined data in pandas.

**GroupBy** A pandas operation that involves splitting data into groups based on some criteria, applying a function to each group, and combining the results.

**Index** An immutable array implementing an ordered, sliceable set for labeling DataFrame and Series objects.

**Boolean Indexing** Using boolean arrays to filter data in pandas objects.

**Vectorization** Performing operations on entire arrays rather than individual elements, for improved performance.

**Method Chaining** Technique of calling multiple methods in sequence on the same object.

**Reshaping** Changing the structure of data from wide to long format or vice versa.

Anshuman Singh — [anshuman365.github.io](https://github.com/anshuman365)

## About the Author



Anshuman Singh

### Anshuman Singh

Data Science Learner & Python Developer

**Name:** Anshuman Singh  
**Email:** anshumansingh3697@gmail.com  
**Website:** anshuman365.github.io  
**GitHub:** github.com/anshuman365  
**LinkedIn:** linkedin.com/in/anshuman-singh-662183303

Anshuman Singh is an enthusiastic data science learner and Python developer on a journey to master the art of data analysis and machine learning. With a passion for turning raw data into meaningful insights, he is constantly exploring new technologies and methodologies in the data science landscape.

Through hands-on projects and continuous learning, Anshuman has developed strong skills in data manipulation with pandas, data visualization, and building machine learning models. He believes in the power of practical learning and enjoys tackling real-world data challenges.

This book represents his learning journey and commitment to sharing knowledge with fellow learners. When not immersed in code and data, he enjoys contributing to open-source projects, writing about his learning experiences, and connecting with the data science community.

*"The beautiful thing about learning is that nobody can take it away from you."*

---

*Keep Learning, Keep Growing  
Connect with me for learning together!*



## Congratulations!

You've completed the Pandas Mastery journey!



You now possess comprehensive pandas skills that will serve you in all your data science endeavors.

Remember, learning is a continuous journey - keep exploring, keep building, and keep growing!

**What's next?** Continue your journey with machine learning, deep learning, and big data technologies!

**Happy Coding!**