# 📖 Data Analysis with Python: From Fundamentals to Advanced Techniques

## Part I: Foundations of Data Analysis

### 1. Introduction to Data Analysis for Python Experts

- What is Data Analysis?
- The Role of Python in Modern Data Analysis
- Understanding the Data Analysis Workflow
- Setting Up a High-Performance Data Analysis Environment

### 2. Core Concepts in Data Analysis

- Understanding Structured, Semi-Structured, and Unstructured Data
- The Data Analysis Pipeline: From Raw Data to Insights
- Exploratory Data Analysis (EDA) vs. Explanatory Analysis
- Data Ethics, Integrity, and Reproducibility

### 3. Python Essentials for Data Analysis *(Quick Recap)*

- Working with Iterators, Generators, and List Comprehensions
- Lambda Functions, Map-Reduce, and Functional Programming Techniques
- Error Handling and Debugging for Data Analysis
- Efficient File I/O: Working with CSV, JSON, and SQL

## Part II: Mastering Numpy for Numerical Computing

### 4. Numpy Essentials

- Numpy Arrays: The Foundation of Numerical Computing
- Vectorization: Replacing Loops with Efficient Operations
- Broadcasting, Indexing, and Slicing Techniques
- Statistical Functions, Aggregations, and Basic Linear Algebra

### 5. Advanced Numpy Techniques

- Memory Management and Data Alignment
- Structured Arrays and Record Arrays
- Optimizing Computations with Numba and Cython
- Random Sampling, Simulations, and Complex Mathematical Operations

## Part III: Data Manipulation with Pandas

### 6. Pandas Fundamentals

- Series, DataFrames, and MultiIndex Structures
- Data Import/Export (CSV, Excel, JSON, SQL, APIs)
- Indexing, Filtering, and Querying Large DataFrames
- Efficient Data Selection and Manipulation

7. **Data Cleaning and Preprocessing**

- Handling Missing, Duplicate, and Inconsistent Data
- Data Transformation: Apply, Map, and Vectorized Functions
- Data Encoding, Normalization, and Type Optimization
- Working with Time Series and Date-Time Indexes

8. **Advanced Data Wrangling**

- GroupBy Operations, Aggregations, and Custom Functions
- Merging, Joining, Concatenating, and Reshaping Data
- Pivot Tables, Cross-Tabulations, and Complex Reshaping
- Handling Nested and Hierarchical Data

9. **Scaling Pandas for Big Data**

- Working with Large Datasets (Out-of-Core Processing)
- Optimizing Memory Usage and Performance
- Using Dask and Vaex for Parallel and Distributed DataFrames
- Chunking and Batch Processing Large Data Files

---

# Part IV: Data Visualization with Matplotlib and Beyond

10. **Fundamentals of Data Visualization**

- Anatomy of a Matplotlib Plot
- Creating Line Plots, Bar Charts, Histograms, and Scatter Plots
- Customizing Visuals: Colors, Labels, Grids, and Annotations
- Subplots, Multiple Axes, and Plot Styling

11. **Advanced Visualization Techniques**

- Seaborn for Statistical Visualizations
- Plotly for Interactive and Web-Based Visualizations
- Heatmaps, Correlation Matrices, and Geospatial Plots
- 3D Plots, Network Graphs, and Complex Visual Representations

12. **Building Interactive Dashboards**

- Integrating Pandas, Matplotlib, and Plotly for Dashboards
- Using ipywidgets and Voila for Interactivity
- Designing Reproducible Reports with Jupyter and nbconvert

---

# Part V: Advanced Techniques for Experts

13. **Optimizing Data Pipelines**

    - Profiling and Benchmarking Data Operations
    - Performance Tuning with Vectorization and Lazy Evaluation
    - Handling Sparse Data and Optimizing Data Types
    - Memory Mapping and Efficient Disk I/O

14. **Memory Management and Scalability**

    - Reducing Memory Footprint in Pandas and Numpy
    - Working with Large-Scale Datasets Using Dask and PySpark
    - Parallel and Distributed Data Analysis Techniques
    - Caching, Streaming Data, and Real-Time Analytics

15. **Advanced Data Wrangling and Transformation**

    - Complex Reshaping Techniques (Melt, Pivot, Stack, Unstack)
    - Advanced GroupBy, Rolling Windows, and Window Functions
    - Time Series Analysis and Anomaly Detection
    - Handling High-Dimensional and Multivariate Data

16. **Dimensionality Reduction and Complex Data Analysis**

    - Principal Component Analysis (PCA) and Singular Value Decomposition (SVD)
    - t-SNE, UMAP, and Other Visualization Techniques
    - Feature Engineering for Predictive Analytics
    - Correlation Analysis and Multicollinearity

# Part VI: Real-World Data Analysis Projects

17. **Exploratory Data Analysis (EDA) in Practice**

    - Designing Effective EDA Workflows
    - Identifying Patterns, Trends, and Anomalies
    - Automating EDA with Custom Scripts
    - Case Study: Analyzing E-Commerce Customer Behavior

18. **Domain-Specific Data Analysis**

    - Financial Data Analysis: Time Series, Forecasting, and Risk Modeling
    - Healthcare Analytics: Patient Data, Statistics, and Predictive Modeling
    - Social Media and Text Data Analytics: Sentiment Analysis with Numpy & Pandas
    - Scientific Computing: Monte Carlo Simulations and Experimental Data Analysis

19. **Building Scalable Data Analysis Pipelines**

    - Structuring Complex Data Projects for Reusability
    - Data Versioning and Reproducibility with DVC and Git
    - Creating Automated ETL Pipelines
    - Case Study: End-to-End Churn Prediction Pipeline

## Part VII: Expanding Your Toolkit

20. **Integrating Data Analysis with External Systems**

    - Working with Databases (SQL, NoSQL) and APIs
    - Cloud Data Processing with AWS, BigQuery, and Snowflake
    - Real-Time Data Streaming with Kafka and Pandas
    - Data Engineering for Data Science Workflows

21. **From Data Analysis to Data Science**

    - Preparing Data for Machine Learning Models
    - Introduction to Scikit-learn for Predictive Analytics
    - Using Pandas and Numpy for Feature Engineering
    - Building and Evaluating Simple Machine Learning Pipelines

22. **Best Practices and Next Steps**

    - Writing Clean, Efficient, and Scalable Code
    - Debugging, Profiling, and Memory Optimization
    - Ensuring Data Quality and Consistency
    - Building Your Data Analysis Portfolio

## 📚 Appendices

- **A. Numpy, Pandas, and Matplotlib Cheat Sheets**
- **B. Performance Optimization Tips for Data Analysts**
- **C. Common Data Analysis Pitfalls and How to Avoid Them**
- **D. Sample Datasets for Practice and Case Studies**
- **E. Resources for Further Learning (Courses, Books, and Tools)**

# 📖 Chapter 1: Introduction to Data Analysis for Python Experts

## 🧠 1. What is Data Analysis?

### ☑️ Definition

Data analysis is the process of inspecting, cleansing, transforming, and modeling data to extract meaningful insights, identify trends, and support decision-making. It bridges the gap between raw data and actionable knowledge.

### ⚡ Key Objectives of Data Analysis:

- **Extract insights** from large, complex datasets
- **Identify patterns, trends, and anomalies**
- **Support data-driven decision-making**
- **Predict future outcomes** using statistical and machine learning models

### 📊 Types of Data Analysis:

| Type | Purpose | Techniques Used |
|---|---|---|
| **Descriptive Analysis** | What happened? | Summarization, Aggregation |
| **Diagnostic Analysis** | Why did it happen? | Correlation, Root Cause Analysis |
| **Predictive Analysis** | What is likely to happen? | Regression, Classification Models |
| **Prescriptive Analysis** | What should we do? | Optimization, Decision Trees |
| **Exploratory Data Analysis (EDA)** | Discover hidden patterns | Visualizations, Clustering |

## 🐍 2. The Role of Python in Modern Data Analysis

### 🏆 Why Python?

Python has become the leading language for data analysis due to its simplicity, extensive ecosystem, and performance optimization techniques.

**Key Strengths:**

- **Readable syntax**: Ideal for rapid prototyping and collaboration.
- **Rich ecosystem**: Libraries like `numpy`, `pandas`, `matplotlib`, and `scikit-learn`.
- **Scalability**: From small datasets to big data frameworks (`Dask`, `PySpark`).
- **Interoperability**: Easy integration with databases, APIs, and cloud services.
- **Community support**: A vast array of tutorials, forums, and open-source projects.

### 💡 Essential Python Libraries for Data Analysis:

| Library | Purpose |
|---|---|
| **Numpy** | High-performance numerical computations |
| **Pandas** | Data manipulation and analysis |
| **Matplotlib** | Basic plotting and visualization |
| **Seaborn** | Statistical graphics built on top of Matplotlib |
| **Scikit-learn** | Machine learning and predictive modeling |
| **Dask** | Parallel computing for large datasets |
| **Plotly** | Interactive visualizations |

## 🚀 Functional Programming in Data Analysis:

Python's functional programming features (like `map`, `filter`, `reduce`, and comprehensions) allow for concise, expressive data transformations.

**Example:** Using `map` and `filter` for data transformation:

```python
# List of sales transactions
sales = [120, 340, 560, 80, 230, 150]

# Apply 10% discount to sales > 200
discounted_sales = list(map(lambda x: x * 0.9 if x > 200 else x, sales))

# Filter out transactions < 100
high_value_sales = list(filter(lambda x: x >= 100, discounted_sales))

print(high_value_sales)
```

**Output:**

```
[108.0, 306.0, 504.0, 230, 150]
```

---

# 🔄 3. Understanding the Data Analysis Workflow

A successful data analysis project follows a systematic workflow to ensure accuracy, scalability, and reproducibility.

## ⚙️ The 6-Step Data Analysis Workflow:

1. **Define Objectives** 🎯

   - Understand the problem and set clear goals.

2. **Collect Data** 🖲️

   - Pull data from databases, APIs, or files (CSV, JSON, Excel).
   - Ensure data is relevant and high-quality.

3. **Clean and Prepare Data** 🧹

   - Handle missing values, duplicates, and data inconsistencies.
   - Normalize, transform, and encode data for analysis.

4. **Explore Data (EDA)** 🔍

   - Use statistical summaries, visualizations, and correlation analysis.
   - Identify trends, patterns, and anomalies.

5. **Analyze and Model** 📊

   - Apply statistical methods or machine learning algorithms.
   - Optimize models for accuracy and performance.

6. **Visualize and Communicate Insights** ☑️

   - Build charts, dashboards, and reports to present findings.
   - Translate data insights into actionable recommendations.

## ⚡ Example: Analyzing Sales Data

```python
import pandas as pd
import matplotlib.pyplot as plt

# Step 1: Load Data
data = pd.DataFrame({
    'Product': ['A', 'B', 'C', 'D', 'E'],
    'Sales': [250, 340, 560, 480, 150],
    'Returns': [5, 7, 9, 3, 4]
})

# Step 2: Calculate Net Sales
data['Net_Sales'] = data['Sales'] - data['Returns']

# Step 3: Visualize
plt.bar(data['Product'], data['Net_Sales'], color='skyblue')
plt.title('Net Sales by Product')
plt.xlabel('Product')
plt.ylabel('Net Sales')
plt.show()
```

**Visualization Output:** A bar chart showing net sales for each product.

---

## 🖥️ 4. Setting Up a High-Performance Data Analysis Environment

## ⚙️ Recommended Tools & Libraries:

- **Jupyter Notebook/Lab** → Interactive coding and visualization
- **VSCode or PyCharm** → Advanced IDEs with data science extensions
- **Anaconda Distribution** → Pre-configured Python data science environment
- **Virtual Environments** → Use venv or conda for dependency management

## 📦 Installing Essential Libraries:

```
pip install numpy pandas matplotlib seaborn scikit-learn jupyter
```

Or using **Anaconda**:

```
conda install numpy pandas matplotlib seaborn scikit-learn jupyter
```

## 🏆 Best Practices for Data Analysis Projects:

- ☑️ **Use version control** (Git) for reproducibility.
- ☑️ **Modularize code** into reusable functions and scripts.
- ☑️ **Document everything** — from data sources to analysis steps.
- ☑️ **Automate repetitive tasks** using functional programming and pipelines.
- ☑️ **Optimize performance** by using vectorized operations and memory-efficient data structures.

## 📌 Chapter Summary:

- **Data analysis** turns raw data into actionable insights.
- **Python** offers an extensive ecosystem for efficient data analysis.
- A **structured workflow** helps ensure clean, scalable, and reproducible analyses.
- **Functional programming techniques** can simplify data manipulation tasks.
- A **high-performance environment** sets the foundation for successful data projects.

# 📖 Chapter 2: Core Concepts in Data Analysis

## 📊 1. Understanding Structured, Semi-Structured, and Unstructured Data

### 🏛 Structured Data

Structured data follows a clear, defined schema, often stored in tabular formats like databases or spreadsheets. It allows for efficient querying and is ideal for traditional data analysis.

- **Examples:**
    - Relational databases (MySQL, PostgreSQL)
    - CSV/Excel files
    - DataFrames in Pandas

**Example in Pandas:**

```python
import pandas as pd

# Structured data in a DataFrame
data = pd.DataFrame({
    'Product': ['A', 'B', 'C'],
    'Price': [100, 150, 200],
    'Stock': [30, 20, 15]
})

print(data)
```

**Output:**

```
   Product  Price  Stock
0        A    100     30
1        B    150     20
2        C    200     15
```

### 🌐 Semi-Structured Data

Semi-structured data doesn't adhere to a strict schema but still contains tags or markers to separate elements.

- **Examples:**
    - JSON, XML
    - NoSQL databases (MongoDB)
    - API responses

**Example (JSON Parsing):**

```python
import json

# Sample JSON data
json_data = '''
{
  "employees": [
    {"name": "Alice", "role": "Developer"},
    {"name": "Bob", "role": "Designer"}
  ]
}
'''

# Parse JSON
parsed = json.loads(json_data)

# Access data
for emp in parsed['employees']:
    print(f"{emp['name']} - {emp['role']}")
```

**Output:**

```
Alice - Developer
Bob - Designer
```

## 🗃 Unstructured Data

Unstructured data lacks a predefined format or organization, making it more challenging to process.

- **Examples:**
    - Text files, emails
    - Images, videos
    - Social media posts

**Text Data Example (Word Count):**

```python
from collections import Counter

text = "Data analysis is fun. Data is powerful."
words = text.lower().replace('.', '').split()
word_count = Counter(words)

print(word_count)
```

**Output:**

```
Counter({'data': 2, 'analysis': 1, 'is': 2, 'fun': 1, 'powerful': 1})
```

---

## 🔄 2. The Data Analysis Pipeline: From Raw Data to Insights

A successful data analysis project follows a logical pipeline, ensuring data integrity and meaningful outcomes.

### 💡 6-Step Data Analysis Pipeline:

1. **Data Collection** 🖱️

   - Gather data from APIs, databases, web scraping, or CSV files.

2. **Data Cleaning** 🖌️

   - Handle missing values, remove duplicates, and correct inconsistencies.

3. **Data Transformation** 🔄

   - Apply scaling, normalization, encoding, and feature engineering.

4. **Exploratory Data Analysis (EDA)** 📊

   - Discover patterns, correlations, and anomalies using visualizations.

5. **Data Modeling & Analysis** ☑️

   - Use statistical methods or machine learning for deeper insights.

6. **Presentation of Results** 📇

   - Share findings using dashboards, reports, or visualizations.

---

### ⚡ Example Pipeline (Sales Data Analysis):

```python
import pandas as pd
import matplotlib.pyplot as plt

# Step 1: Load Data
data = pd.DataFrame({
    'Product': ['A', 'B', 'C', 'D'],
    'Sales': [250, 340, 560, 480],
    'Returns': [5, 7, 9, 3]
})

# Step 2: Clean Data (No missing values here)

# Step 3: Transform Data (Net Sales)
data['Net_Sales'] = data['Sales'] - data['Returns']
```

```python
# Step 4: EDA (Visualize Net Sales)
plt.bar(data['Product'], data['Net_Sales'], color='lightgreen')
plt.title('Net Sales by Product')
plt.xlabel('Product')
plt.ylabel('Net Sales')
plt.show()
```

# 🔍 3. Exploratory Data Analysis (EDA) vs. Explanatory Analysis

## ⚛ Exploratory Data Analysis (EDA):

- Purpose: **Discover patterns, spot anomalies, and test hypotheses.**
- Methods: Descriptive statistics, visualizations, correlation matrices.
- Outcome: Identify trends and areas for deeper investigation.

**Example (EDA with Pandas):**

```python
import seaborn as sns

# Sample Data
df = sns.load_dataset('tips')

# Correlation heatmap
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
plt.title('Correlation Matrix')
plt.show()
```

## 📖 Explanatory Data Analysis:

- Purpose: **Communicate specific insights and findings.**
- Methods: Polished visualizations, storytelling with data.
- Outcome: Clear and concise presentation of results to stakeholders.

**Key Differences:**

| Aspect | EDA | Explanatory Analysis |
|---|---|---|
| **Goal** | Explore, discover patterns | Communicate insights |
| **Audience** | Analysts, Data Scientists | Stakeholders, Decision-Makers |
| **Tools Used** | Histograms, Boxplots, Heatmaps | Dashboards, Infographics |
| **Flexibility** | Open-ended | Goal-driven |

# ⚖️ 4. Data Ethics, Integrity, and Reproducibility

## 📋 Data Ethics:

- **Transparency:** Always disclose data sources and methods.
- **Privacy:** Follow data protection laws (e.g., GDPR).
- **Fairness:** Avoid biases that could impact results.
- **Accountability:** Ensure responsible data handling.

---

## 📊 Data Integrity:

- **Data Accuracy:** Ensure data reflects real-world conditions.
- **Consistency:** Maintain uniform formats and standards.
- **Completeness:** Fill in missing data or acknowledge gaps.

**Example: Handling Missing Data in Pandas:**

```python
import pandas as pd

# Sample DataFrame with missing values
df = pd.DataFrame({
    'Name': ['Alice', 'Bob', 'Charlie'],
    'Score': [85, None, 92]
})

# Fill missing values with the mean
df['Score'].fillna(df['Score'].mean(), inplace=True)

print(df)
```

**Output:**

```
      Name  Score
0    Alice   85.0
1      Bob   88.5
2  Charlie   92.0
```

---

## 🔁 Reproducibility in Data Analysis:

- **Version Control:** Use `Git` to track changes in data and code.
- **Documentation:** Keep clear records of data sources, transformations, and decisions.
- **Automated Pipelines:** Build repeatable workflows using scripts or notebooks.
- **Notebook Best Practices:** Use Jupyter with clear markdown cells, code explanations, and outputs.

**Example: Reproducible Analysis Workflow**

1. **Data Collection:** API calls → Save raw data
2. **Data Cleaning:** Scripted transformations → Save cleaned dataset
3. **Modeling & Analysis:** Version-controlled notebooks
4. **Results Sharing:** Export plots/reports → Share as PDFs or interactive dashboards

## 📌 Chapter Summary:

- **Structured, semi-structured, and unstructured data** require different analysis techniques.
- The **Data Analysis Pipeline** transforms raw data into insights through cleaning, transformation, and modeling.
- **EDA** helps discover patterns, while **Explanatory Analysis** communicates results effectively.
- Upholding **data ethics, integrity, and reproducibility** ensures trust in data-driven decisions.

# 📖 Chapter 3: Python Essentials for Data Analysis
## *(Quick Recap)*

This chapter provides a concise yet powerful recap of essential Python concepts tailored for data analysis. While you're already a Python expert, this chapter will focus on applying core Python techniques effectively in data analysis workflows.

## 🔁 1. Working with Iterators, Generators, and List Comprehensions

### 🎯 Iterators

An **iterator** is an object that enables traversing through a collection, one element at a time, without needing to store the entire collection in memory.

- **Key methods:**
  - `__iter__()` — returns the iterator object.
  - `__next__()` — returns the next element.

**Example: Custom Iterator**

```python
class Counter:
    def __init__(self, low, high):
        self.current = low
        self.high = high

    def __iter__(self):
        return self

    def __next__(self):
        if self.current > self.high:
            raise StopIteration
        else:
            self.current += 1
            return self.current - 1

# Using the iterator
for num in Counter(1, 5):
    print(num)
```

**Output:**

```
1
2
3
4
5
```

## ⚡ Generators

Generators simplify iterators using the `yield` statement. They're memory-efficient, ideal for large datasets.

**Example: Data Stream Generator**

```python
def data_stream(n):
    for i in range(n):
        yield i * 2

stream = data_stream(5)
for value in stream:
    print(value)
```

**Output:**

```
0
2
4
6
8
```

💡 **Use case in Data Analysis:**

- Streaming large datasets row by row to avoid memory overload.

## 🧮 List Comprehensions

List comprehensions offer a concise way to create lists and apply transformations.

**Example: Filtering Data**

```python
# Filter even numbers from a list
numbers = [1, 2, 3, 4, 5, 6]
evens = [x for x in numbers if x % 2 == 0]
print(evens)
```

**Output:**

```
[2, 4, 6]
```

💡 **Tip:** Use **set comprehensions** `{}` and **dict comprehensions** `{key: value}` similarly.

---

# 🏛 2. Lambda Functions, Map-Reduce, and Functional Programming Techniques

## 🌀 Lambda Functions

Lambda functions are anonymous, inline functions, perfect for short transformations.

**Example: Simple Lambda**

```python
square = lambda x: x ** 2
print(square(5))
```

**Output:**

```
25
```

---

## 🗺 Map, Filter, and Reduce

- `map()` → Applies a function to each item in a list.
- `filter()` → Filters items based on a condition.
- `reduce()` → Aggregates values (requires `functools`).

**Example: Apply Map and Filter**

```python
from functools import reduce

numbers = [1, 2, 3, 4, 5]

# Double each number
doubled = list(map(lambda x: x * 2, numbers))

# Filter even numbers
evens = list(filter(lambda x: x % 2 == 0, doubled))

# Sum all numbers
total = reduce(lambda x, y: x + y, evens)

print(f"Doubled: {doubled}")
print(f"Evens: {evens}")
print(f"Sum of Evens: {total}")
```

**Output:**

```
Doubled: [2, 4, 6, 8, 10]
Evens: [2, 4, 6, 8, 10]
Sum of Evens: 30
```

## 💡 Functional Programming in Data Analysis

Functional programming techniques simplify data pipelines by applying operations in sequence.

**Example: Chaining Operations with `map` and `filter`**

```python
data = range(10)

result = sum(
    map(lambda x: x ** 2,
        filter(lambda x: x % 2 == 0, data))
)

print(result)
```

**Explanation:**

1. **Filter** even numbers.
2. **Square** the filtered numbers.
3. **Sum** the squared values.

**Output:**

```
120
```

## 🚂 3. Error Handling and Debugging for Data Analysis

Working with real-world data means encountering errors — missing values, malformed data, or unexpected types. Efficient error handling is crucial.

### ⚠ Try-Except Blocks

Use `try-except` to catch runtime errors gracefully.

**Example: Handling Division by Zero**

```python
def safe_divide(a, b):
    try:
        return a / b
    except ZeroDivisionError:
```

```
        return "Cannot divide by zero"

print(safe_divide(10, 2))  # 5.0
print(safe_divide(10, 0))  # Cannot divide by zero
```

## 🔍 Debugging Tips for Data Analysis:

- **Use assertions** to catch anomalies early:

  ```
  assert df.isnull().sum().sum() == 0, "Data contains missing values!"
  ```

- **Leverage pdb for step-by-step debugging:**

  ```
  import pdb
  pdb.set_trace()
  ```

- **Use logging instead of print statements for better traceability:**

  ```
  import logging

  logging.basicConfig(level=logging.INFO)
  logging.info("Data loaded successfully")
  ```

# 📁 4. Efficient File I/O: Working with CSV, JSON, and SQL

## 📄 CSV Files

**Read and write CSV using Pandas:**

```
import pandas as pd

# Read CSV
df = pd.read_csv('data.csv')

# Write CSV
df.to_csv('output.csv', index=False)
```

## 📊 JSON Files

**Load and save JSON data:**

```python
import json

# Sample data
data = {"name": "Alice", "age": 30, "role": "Analyst"}

# Write to JSON
with open('data.json', 'w') as file:
    json.dump(data, file)

# Read JSON
with open('data.json', 'r') as file:
    loaded_data = json.load(file)

print(loaded_data)
```

## 🗄 SQL Databases

**Connect and query using `sqlite3`:**

```python
import sqlite3
import pandas as pd

# Create connection
conn = sqlite3.connect('example.db')

# Write DataFrame to SQL
df = pd.DataFrame({
    'Name': ['Alice', 'Bob'],
    'Age': [30, 25]
})
df.to_sql('users', conn, if_exists='replace', index=False)

# Read from SQL
query = "SELECT * FROM users"
result = pd.read_sql(query, conn)
print(result)

# Close connection
conn.close()
```

## 📌 Chapter Summary:

- **Iterators, Generators, and List Comprehensions** streamline data workflows.
- **Lambda, Map-Reduce, and Functional Programming** simplify complex data transformations.
- **Robust error handling** ensures data integrity during analysis.
- Efficient **File I/O operations** enable seamless data movement between formats and databases.

# 📖 Chapter 4: Numpy Essentials

NumPy is the cornerstone of numerical computing in Python. Its efficient array structures and high-performance operations make it a foundational tool for data analysis, scientific computing, and machine learning. This chapter will cover essential NumPy concepts and advanced techniques to help you leverage its full potential.

## 📐 1. Numpy Arrays: The Foundation of Numerical Computing

### ☑ Why NumPy Arrays Over Python Lists?

- **Memory Efficiency:** NumPy arrays use less memory compared to Python lists.
- **Speed:** Operations on arrays are significantly faster due to optimized C backend.
- **Convenience:** Built-in functions for mathematical operations, broadcasting, and aggregation.

### 🛠 Creating Numpy Arrays

```python
import numpy as np

# 1D Array
arr1 = np.array([1, 2, 3, 4])
print("1D Array:", arr1)

# 2D Array (Matrix)
arr2 = np.array([[1, 2], [3, 4]])
print("2D Array:\n", arr2)

# Special Arrays
zeros = np.zeros((3, 3))        # 3x3 matrix of zeros
ones = np.ones((2, 4))          # 2x4 matrix of ones
identity = np.eye(3)            # 3x3 identity matrix
random_arr = np.random.rand(2, 3) # 2x3 random numbers

print("Zeros:\n", zeros)
print("Random Array:\n", random_arr)
```

### 📊 Key Properties of Numpy Arrays

```python
print("Shape:", arr2.shape)      # (2, 2)
print("Size:", arr2.size)        # 4
print("Data Type:", arr2.dtype)  # int64 or float64
```

## ⚡ 2. Vectorization: Replacing Loops with Efficient Operations

**Vectorization** allows operations on entire arrays without explicit loops, resulting in concise and faster code.

## VS For Loop vs. Vectorized Operation

```python
# Python loop
data = [1, 2, 3, 4]
squared_loop = [x**2 for x in data]

# NumPy vectorized
arr = np.array(data)
squared_np = arr ** 2

print("Loop Result:", squared_loop)
print("NumPy Result:", squared_np)
```

## ⏱ Performance Comparison:

```python
import time

data = np.random.rand(1000000)

# Python list
start = time.time()
squared_list = [x**2 for x in data]
print("Python list time:", time.time() - start)

# NumPy array
arr = np.array(data)
start = time.time()
squared_np = arr ** 2
print("NumPy time:", time.time() - start)
```

💡 **NumPy will be ~10-50x faster** due to optimized C routines.

---

## ✨ Common Vectorized Operations

```python
arr = np.array([10, 20, 30, 40])

# Arithmetic
print(arr + 5)        # [15 25 35 45]
print(arr * 2)        # [20 40 60 80]

# Element-wise operations
print(np.sqrt(arr))   # Square root
print(np.log(arr))    # Natural log
```

```python
# Conditional operations
print(arr[arr > 20])  # [30 40]
```

---

# 📏 3. Broadcasting, Indexing, and Slicing Techniques

## 📡 Broadcasting

Broadcasting allows NumPy to perform arithmetic operations on arrays of different shapes.

```python
arr = np.array([[1, 2, 3],
                [4, 5, 6]])

# Add a scalar to entire array
print(arr + 10)

# Broadcasting with arrays
col_vector = np.array([[1], [2]])
result = arr + col_vector

print("Broadcasted Result:\n", result)
```

### 💡 Rules of Broadcasting:

1. Dimensions are compared from right to left.
2. A dimension of size 1 can be stretched to match.

---

## 🖼 Indexing and Slicing

```python
arr = np.array([[10, 20, 30],
                [40, 50, 60],
                [70, 80, 90]])

# Basic Slicing
print(arr[1, :])        # Row 1 → [40 50 60]
print(arr[:, 2])        # Column 2 → [30 60 90]

# Boolean Indexing
print(arr[arr > 50])    # [60 70 80 90]

# Fancy Indexing
rows = np.array([0, 2])
cols = np.array([1, 2])
print(arr[rows, cols])  # [20 90]
```

---

## 🐾 Advanced Slicing Techniques

- **Reversing Arrays:**

```python
arr = np.arange(10)
print(arr[::-1])  # [9 8 7 6 5 4 3 2 1 0]
```

- **Selecting Diagonals:**

```python
matrix = np.arange(1, 10).reshape(3, 3)
print(np.diag(matrix))  # [1 5 9]
```

# 📊 4. Statistical Functions, Aggregations, and Basic Linear Algebra

## ☑ Statistical Functions

```python
data = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9])

print("Mean:", np.mean(data))
print("Median:", np.median(data))
print("Standard Deviation:", np.std(data))
print("Sum:", np.sum(data))
print("Min:", np.min(data))
print("Max:", np.max(data))
```

## 📊 Aggregations Along Axes

```python
matrix = np.array([[1, 2, 3],
                   [4, 5, 6],
                   [7, 8, 9]])

print("Sum along rows:", np.sum(matrix, axis=1))  # [ 6 15 24]
print("Sum along columns:", np.sum(matrix, axis=0)) # [12 15 18]
```

## 📐 Basic Linear Algebra with NumPy

NumPy provides efficient implementations for common linear algebra operations.

```python
from numpy.linalg import inv, eig, det

# Matrix operations
A = np.array([[1, 2], [3, 4]])
```

```python
# Transpose
print("Transpose:\n", A.T)

# Determinant
print("Determinant:", det(A))

# Inverse
print("Inverse:\n", inv(A))

# Eigenvalues and Eigenvectors
values, vectors = eig(A)
print("Eigenvalues:", values)
print("Eigenvectors:\n", vectors)
```

## 🧪 Solving Linear Systems

```python
# Solve Ax = b
A = np.array([[3, 1], [1, 2]])
b = np.array([9, 8])

x = np.linalg.solve(A, b)
print("Solution x:", x)
```

# 📌 Chapter Summary:

- **NumPy arrays** provide efficient storage and computation for numerical data.
- **Vectorization** eliminates the need for slow Python loops.
- **Broadcasting** enables operations on arrays of different shapes.
- Use **indexing, slicing, and boolean masks** for flexible data access.
- Leverage **statistical functions and linear algebra operations** for complex analysis.

# 📖 Chapter 5: Advanced Numpy Techniques

In this chapter, we'll go beyond the basics of NumPy and dive into advanced topics that unlock the full power of numerical computing in Python. We will explore memory management strategies, structured arrays, optimization techniques using **Numba** and **Cython**, and how to perform simulations and complex mathematical operations.

## 🧠 1. Memory Management and Data Alignment

### 🏛 Understanding NumPy's Memory Model

NumPy arrays are stored in contiguous blocks of memory, which enables fast access and manipulation. Efficient memory usage is crucial for handling large datasets.

### 📊 Key Concepts:

- **Contiguity:** NumPy arrays can be *C-contiguous* (row-major) or *F-contiguous* (column-major).
- **Views vs. Copies:**
    - **Views** share memory with the original array (faster, less memory).
    - **Copies** allocate new memory (safe for independent modifications).

**Example: View vs. Copy**

```python
import numpy as np

arr = np.array([1, 2, 3, 4, 5])
view = arr[1:4]   # This is a view
view[0] = 99
print("Original Array:", arr)  # [1 99 3 4 5]

copy = arr[1:4].copy()
copy[0] = 42
print("After Copy Modification:", arr)  # [1 99 3 4 5]
```

### 🗄 Optimizing Memory Usage

- **Data Types (`dtype`) Optimization:**

```python
# Using smaller data types
arr_float64 = np.arange(1e6, dtype=np.float64)
arr_float32 = np.arange(1e6, dtype=np.float32)

print("Float64 size:", arr_float64.nbytes, "bytes")
print("Float32 size:", arr_float32.nbytes, "bytes")
```

- **Using** `np.memmap` **for Large Datasets:** Memory-map large files to avoid loading everything into RAM.

```python
data = np.memmap('large_array.dat', dtype='float32', mode='w+', shape=(10000,
10000))
data[5000, 5000] = 42
data.flush()  # Ensure data is written to disk
```

---

## 🗄 2. Structured Arrays and Record Arrays

Structured arrays allow you to store complex datasets with mixed data types — similar to SQL tables or Pandas DataFrames but with lower-level control.

### 🛠 Creating Structured Arrays:

```python
# Define a structured dtype
dt = np.dtype([('Name', 'U10'), ('Age', 'i4'), ('Salary', 'f4')])

# Create the array
data = np.array([('Alice', 25, 50000.0), ('Bob', 30, 60000.0)], dtype=dt)

print(data)
```

**Output:**

```
[('Alice', 25, 50000.) ('Bob', 30, 60000.)]
```

### 🔍 Accessing Structured Data:

```python
# Access specific fields
print(data['Name'])   # ['Alice' 'Bob']
print(data['Salary']) # [50000. 60000.]

# Conditional filtering
print(data[data['Age'] > 26])  # [('Bob', 30, 60000.)]
```

---

### 🔁 Record Arrays (`recarray`) for Attribute Access:

```python
rec_data = data.view(np.recarray)
print(rec_data.Name)  # ['Alice' 'Bob']
```

💡 **When to use:**

- Use structured arrays for lightweight data manipulation.
- Use Pandas for heavy-duty data analysis and complex indexing.

---

# ⚡ 3. Optimizing Computations with Numba and Cython

When pure NumPy isn't fast enough, tools like **Numba** and **Cython** can supercharge performance.

## 🚀 Numba: Just-In-Time Compilation

**Installation:**

```
pip install numba
```

**Example: Speeding Up Loops**

```python
from numba import jit
import numpy as np
import time

@jit(nopython=True)
def slow_sum(arr):
    total = 0
    for x in arr:
        total += x
    return total

arr = np.random.rand(10**7)

# Without Numba
start = time.time()
print(np.sum(arr))
print("Numpy Time:", time.time() - start)

# With Numba
start = time.time()
print(slow_sum(arr))
print("Numba Time:", time.time() - start)
```

💡 **Result:** Numba often achieves **5-50x** speedups over regular Python loops.

---

## 🐍 Cython: C Extensions for Python

**Installation:**

```
pip install cython
```

**Example: Accelerating Python with Cython**

1. Create a `cython_example.pyx` file:

```
def cython_sum(double[:] arr):
    cdef double total = 0
    for i in range(arr.shape[0]):
        total += arr[i]
    return total
```

2. Compile with:

```
cythonize -i cython_example.pyx
```

3. Use in Python:

```python
import numpy as np
from cython_example import cython_sum

arr = np.random.rand(10**7)
print(cython_sum(arr))
```

💡 **Cython** is great for loops that can't be vectorized easily.

---

# 🎲 4. Random Sampling, Simulations, and Complex Mathematical Operations

## 🎰 Random Sampling with `numpy.random`

```python
rng = np.random.default_rng(seed=42)

# Random integers
ints = rng.integers(low=0, high=10, size=5)
print("Random Integers:", ints)

# Random floats
floats = rng.random(size=5)
print("Random Floats:", floats)

# Normal distribution
```

```python
normal = rng.normal(loc=0, scale=1, size=5)
print("Normal Distribution:", normal)
```

---

## 🧪 Monte Carlo Simulation Example

Estimate π using random sampling:

```python
import numpy as np

def estimate_pi(n_samples=1000000):
    x = np.random.uniform(0, 1, n_samples)
    y = np.random.uniform(0, 1, n_samples)
    inside_circle = (x**2 + y**2) <= 1
    pi_estimate = 4 * np.sum(inside_circle) / n_samples
    return pi_estimate

print("Estimated π:", estimate_pi())
```

---

## 🔢 Complex Mathematical Operations

NumPy supports advanced math operations directly:

- **Fourier Transforms:**

  ```python
  from numpy.fft import fft

  x = np.random.rand(100)
  y = fft(x)
  print("FFT Result:", y)
  ```

- **Polynomials:**

  ```python
  p = np.poly1d([1, 2, 1])  # x^2 + 2x + 1
  print(p(3))  # Evaluate at x=3
  ```

- **Matrix Factorization (SVD):**

  ```python
  A = np.random.rand(3, 3)
  U, S, Vt = np.linalg.svd(A)
  print("Singular Values:", S)
  ```

## 📌 Chapter Summary:

- **Memory management** helps optimize large datasets.
- **Structured arrays** provide flexible data storage for mixed types.
- **Numba and Cython** accelerate computations for heavy-duty tasks.
- **Random sampling and simulations** enable probabilistic modeling and statistical analysis.

# 📖 Chapter 6: Pandas Fundamentals

**Pandas** is the go-to library for data manipulation and analysis in Python. Its core data structures—**Series**, **DataFrames**, and **MultiIndex**—make working with structured data simple and intuitive. In this chapter, we'll cover the fundamental components of Pandas and how to efficiently import, query, and manipulate large datasets.

## 📦 1. Series, DataFrames, and MultiIndex Structures

### 📊 1.1 Series: The One-Dimensional Labeled Array

A **Series** is a one-dimensional array with labels (indexes).

```python
import pandas as pd

# Creating a Series
data = pd.Series([10, 20, 30, 40], index=['a', 'b', 'c', 'd'])
print(data)

# Accessing by label
print("Element 'b':", data['b'])

# Vectorized operations
print("Doubled values:\n", data * 2)
```

**Output:**

```
a    10
b    20
c    30
d    40
dtype: int64
Element 'b': 20
Doubled values:
a    20
b    40
c    60
d    80
dtype: int64
```

💡 **Tip:** Series can be treated like both a NumPy array (for math ops) and a dictionary (for key-based access).

### 📋 1.2 DataFrames: The 2D Data Powerhouse

A **DataFrame** is a 2D labeled data structure with rows and columns.

```python
# Creating a DataFrame
data = {'Name': ['Alice', 'Bob', 'Charlie'],
        'Age': [25, 30, 35],
        'Salary': [70000, 80000, 90000]}

df = pd.DataFrame(data)
print(df)
```

**Output:**

```
      Name  Age  Salary
0    Alice   25   70000
1      Bob   30   80000
2  Charlie   35   90000
```

## 🔑 Key DataFrame Operations:

```python
# Access columns
print(df['Name'])

# Add a new column
df['Bonus'] = df['Salary'] * 0.10

# Basic statistics
print(df.describe())

# Renaming columns
df.rename(columns={'Salary': 'Annual Salary'}, inplace=True)
print(df)
```

## 🏷️ 1.3 MultiIndex: Working with Hierarchical Data

**MultiIndex** enables multi-level (hierarchical) indexing in DataFrames.

```python
# Creating a MultiIndex DataFrame
arrays = [['USA', 'USA', 'Canada', 'Canada'], ['NY', 'CA', 'ON', 'QC']]
index = pd.MultiIndex.from_arrays(arrays, names=('Country', 'State'))

data = pd.DataFrame({'Population': [20, 40, 10, 15]}, index=index)
print(data)
```

**Output:**

```
                    Population
Country State
USA      NY                20
         CA                40
Canada   ON                10
         QC                15
```

**Accessing with MultiIndex:**

```python
print(data.loc['USA'])
```

---

## 🗂 2. Data Import/Export (CSV, Excel, JSON, SQL, APIs)

Pandas provides powerful I/O tools for reading/writing data from various sources.

### 📄 2.1 CSV Files

```python
# Read CSV
df = pd.read_csv('data.csv')

# Write to CSV
df.to_csv('output.csv', index=False)
```

### 📊 2.2 Excel Files

```python
# Read Excel
df = pd.read_excel('data.xlsx', sheet_name='Sheet1')

# Write to Excel
df.to_excel('output.xlsx', index=False)
```

### 🌐 2.3 JSON Files

```python
# Read JSON
df = pd.read_json('data.json')

# Convert DataFrame to JSON
df.to_json('output.json', orient='records')
```

### 🗄 2.4 SQL Databases

```python
import sqlite3

# Connect to SQLite database
conn = sqlite3.connect('my_database.db')

# Read SQL table
df = pd.read_sql('SELECT * FROM employees', conn)

# Write DataFrame to SQL
df.to_sql('new_table', conn, if_exists='replace', index=False)
```

## 📡 2.5 APIs (Using JSON)

```python
import requests

response = requests.get('https://jsonplaceholder.typicode.com/posts')
json_data = response.json()

# Convert API data to DataFrame
df = pd.DataFrame(json_data)
print(df.head())
```

# 🔍 3. Indexing, Filtering, and Querying Large DataFrames

## 📌 3.1 Indexing Techniques

```python
# Setting index
df.set_index('Name', inplace=True)

# Resetting index
df.reset_index(inplace=True)
```

## 📋 3.2 Filtering Data

```python
# Boolean indexing
filtered_df = df[df['Age'] > 30]

# Multiple conditions
filtered_df = df[(df['Age'] > 25) & (df['Salary'] > 75000)]
```

## 💡 3.3 Efficient Querying with `.query()`

```python
# Using query for better readability
high_salary = df.query('Salary > 75000')
print(high_salary)
```

## 🚀 3.4 Handling Large Datasets

- **Load data in chunks:**

```python
chunks = pd.read_csv('large_data.csv', chunksize=100000)
for chunk in chunks:
    print(chunk.shape)
```

- **Optimize memory usage:**

```python
df = pd.read_csv('large_data.csv', dtype={'ID': 'int32', 'Salary':
'float32'})
print(df.info())
```

# ✂️ 4. Efficient Data Selection and Manipulation

## 📖 4.1 Selecting Data with `.loc[]` and `.iloc[]`

- `loc[]` → **label-based selection**
- `iloc[]` → **position-based selection**

```python
# loc example
print(df.loc[0, 'Name'])

# iloc example
print(df.iloc[1, 2])  # Row 1, Column 2
```

## 🛠️ 4.2 Data Transformation

- **Apply functions to columns:**

```python
df['Age Plus 5'] = df['Age'].apply(lambda x: x + 5)
```

- **Vectorized string operations:**

```python
df['Name'] = df['Name'].str.upper()
```

- **Handling missing data:**

```python
df.fillna(0, inplace=True)        # Replace NaN with 0
df.dropna(inplace=True)           # Drop rows with NaN
```

## 📊 4.3 Grouping and Aggregating Data

```python
# Group by and aggregate
grouped = df.groupby('Department')['Salary'].agg(['mean', 'max', 'min'])
print(grouped)
```

## 🗓 4.4 Pivot Tables and Cross Tabulations

```python
# Pivot Table
pivot = pd.pivot_table(df, values='Salary', index='Department', columns='Gender', aggfunc='mean')

# Cross Tabulation
cross_tab = pd.crosstab(df['Department'], df['Gender'])
```

# 📌 Chapter Summary:

- **Series** and **DataFrames** form the foundation of Pandas data structures.
- Efficiently **import/export** data from various sources (CSV, Excel, JSON, SQL, APIs).
- Mastered **indexing, filtering, and querying** for small and large datasets.
- Learned **grouping, aggregation**, and advanced data manipulation techniques.

# 📖 Chapter 7: Data Cleaning and Preprocessing

**Data cleaning and preprocessing** is a crucial step in any data analysis pipeline. Dirty data can lead to inaccurate insights and misleading conclusions. In this chapter, we'll cover strategies to handle missing, duplicate, and inconsistent data, apply data transformations, encode and normalize data, optimize data types, and work effectively with time series data.

## 🖌️ 1. Handling Missing, Duplicate, and Inconsistent Data

### ❓ 1.1 Identifying Missing Data

Pandas uses **NaN** to represent missing data.

```python
import pandas as pd
import numpy as np

# Sample DataFrame
df = pd.DataFrame({
    'Name': ['Alice', 'Bob', 'Charlie', 'David', np.nan],
    'Age': [25, 30, np.nan, 40, 35],
    'Salary': [70000, 80000, 90000, np.nan, 60000]
})

# Detect missing values
print(df.isnull())

# Count missing values per column
print(df.isnull().sum())
```

**Output:**

```
    Name    Age  Salary
0  False  False   False
1  False  False   False
2  False   True   False
3  False  False    True
4   True  False   False

Name    1
Age     1
Salary  1
dtype: int64
```

### 🛠️ 1.2 Handling Missing Data

- **Drop missing values:**

```python
# Drop rows with any missing values
df_drop = df.dropna()

# Drop columns with missing values
df_drop_col = df.dropna(axis=1)
```

- **Fill missing values:**

```python
# Fill with a constant
df_fill = df.fillna(0)

# Forward fill (propagate last valid observation)
df_ffill = df.fillna(method='ffill')

# Backward fill
df_bfill = df.fillna(method='bfill')

# Fill with mean
df['Age'].fillna(df['Age'].mean(), inplace=True)
```

## ⬜ 1.3 Handling Duplicates

```python
# Sample DataFrame with duplicates
df_dup = pd.DataFrame({
    'Name': ['Alice', 'Bob', 'Alice', 'David', 'Bob'],
    'Age': [25, 30, 25, 40, 30]
})

# Detect duplicates
print(df_dup.duplicated())

# Drop duplicate rows
df_unique = df_dup.drop_duplicates()
```

## ⚖️ 1.4 Resolving Inconsistent Data

- **Fix inconsistent casing:**

```python
df['Name'] = df['Name'].str.title()  # 'alice' -> 'Alice'
```

- **Remove extra spaces:**

```python
df['Name'] = df['Name'].str.strip()
```

- **Replace inconsistent entries:**

```python
df['Name'].replace({'Alicia': 'Alice'}, inplace=True)
```

# 🔄 2. Data Transformation: Apply, Map, and Vectorized Functions

### 🧮 2.1 Vectorized Operations

Pandas leverages NumPy under the hood, enabling **vectorized operations** for speed.

```python
df['Salary After Tax'] = df['Salary'] * 0.7
```

### 🗺️ 2.2 Map, Apply, and Applymap

- `map()` → element-wise operation for Series
- `apply()` → apply function across axis (rows/columns)
- `applymap()` → element-wise operation for entire DataFrame

```python
# Using map on a Series
df['Age Category'] = df['Age'].map(lambda x: 'Young' if x < 30 else 'Old')

# Using apply on DataFrame
df['Tax'] = df['Salary'].apply(lambda x: x * 0.3 if pd.notnull(x) else 0)

# applymap for entire DataFrame
df_numeric = df[['Age', 'Salary']]
df_scaled = df_numeric.applymap(lambda x: x / 1000)
```

### 🔬 2.3 Complex Transformations with `pipe()`

`pipe()` enables functional-style chaining for cleaner code.

```python
def add_bonus(df, bonus):
    df['Salary'] += bonus
    return df

df = df.pipe(add_bonus, bonus=5000)
```

# 🔢 3. Data Encoding, Normalization, and Type Optimization

## 📋 3.1 Encoding Categorical Variables

- **One-Hot Encoding:**

```python
df = pd.DataFrame({'Gender': ['Male', 'Female', 'Male']})
df_encoded = pd.get_dummies(df, columns=['Gender'], drop_first=True)
```

- **Label Encoding:**

```python
df['Gender_Code'] = df['Gender'].map({'Male': 1, 'Female': 0})
```

## 📏 3.2 Normalization and Scaling

- **Min-Max Scaling:**

```python
from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()
df[['Salary']] = scaler.fit_transform(df[['Salary']])
```

- **Standardization (Z-score):**

```python
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
df[['Age']] = scaler.fit_transform(df[['Age']])
```

## 💾 3.3 Optimizing Data Types for Memory Efficiency

```python
# Check memory usage
print(df.info(memory_usage='deep'))

# Convert float64 to float32
df['Salary'] = df['Salary'].astype('float32')

# Convert object to category
df['Gender'] = df['Gender'].astype('category')
```

💡 **Tip:** Optimizing data types can significantly reduce memory usage when working with large datasets.

---

# 📅 4. Working with Time Series and Date-Time Indexes

## ⏰ 4.1 Parsing and Converting Dates

```python
# Convert to datetime
df['Join Date'] = pd.to_datetime(df['Join Date'])

# Extract date parts
df['Year'] = df['Join Date'].dt.year
df['Month'] = df['Join Date'].dt.month
```

---

## 📅 4.2 Setting Date-Time as Index

```python
# Set date as index for time series analysis
df.set_index('Join Date', inplace=True)

# Resample to monthly data
df_monthly = df.resample('M').mean()
```

---

## 📊 4.3 Time Series Operations

- **Date Range Creation:**

```python
# Create a range of dates
date_rng = pd.date_range(start='2023-01-01', end='2023-12-31', freq='M')
```

- **Rolling Window Analysis:**

```python
# 3-month rolling average
df['Rolling_Mean'] = df['Salary'].rolling(window=3).mean()
```

- **Shifting Data:**

```python
# Lagging by 1 period
df['Salary_Lag1'] = df['Salary'].shift(1)
```

---

# 📌 Chapter Summary:

- Handled **missing, duplicate, and inconsistent data** using Pandas functions.
- Applied **data transformations** using `map()`, `apply()`, and vectorized operations.
- Performed **data encoding, normalization**, and optimized data types for efficiency.
- Managed **time series data** with datetime indexes and rolling statistics.

---

# 📖 Chapter 8: Advanced Data Wrangling

Data wrangling is the process of transforming raw data into a structured and clean format for analysis. In this chapter, we'll dive deep into **grouping**, **aggregating**, **merging**, **reshaping**, and **handling complex data structures** using Pandas. You'll learn how to efficiently manipulate large datasets, enabling you to extract meaningful insights.

---

## 📊 1. GroupBy Operations, Aggregations, and Custom Functions

The `groupby()` method in Pandas is a powerful tool for splitting data into groups based on specific criteria, performing operations on each group, and then combining the results.

### 📂 1.1 Basic GroupBy and Aggregation

```python
import pandas as pd

# Sample DataFrame
df = pd.DataFrame({
    'Department': ['HR', 'HR', 'IT', 'IT', 'Finance', 'Finance'],
    'Employee': ['Alice', 'Bob', 'Charlie', 'David', 'Eva', 'Frank'],
    'Salary': [70000, 80000, 90000, 85000, 75000, 72000],
    'Bonus': [5000, 6000, 8000, 7500, 4000, 4500]
})

# Group by Department and calculate mean salary
grouped = df.groupby('Department')['Salary'].mean()
print(grouped)
```

**Output:**

```
Department
Finance    73500.0
HR         75000.0
IT         87500.0
Name: Salary, dtype: float64
```

---

### 🔢 1.2 Multiple Aggregations with `agg()`

```python
# Apply multiple aggregations
agg_result = df.groupby('Department').agg({
    'Salary': ['mean', 'max', 'min'],
    'Bonus': 'sum'
```

```
})
print(agg_result)
```

**Output:**

```
           Salary                        Bonus
             mean     max     min     sum
Department
Finance    73500.0  75000  72000    8500
HR         75000.0  80000  70000   11000
IT         87500.0  90000  85000   15500
```

## 🛠️ 1.3 Custom Aggregations with `apply()`

You can use `apply()` to run custom functions on grouped data.

```python
# Custom function to calculate range
def salary_range(x):
    return x.max() - x.min()

# Apply custom function
df_grouped = df.groupby('Department')['Salary'].apply(salary_range)
print(df_grouped)
```

**Output:**

```
Department
Finance     3000
HR         10000
IT          5000
Name: Salary, dtype: int64
```

## 🧮 1.4 Transforming Data with `transform()`

- `transform()` returns a DataFrame of the same shape as the input, allowing for broadcasting operations.

```python
# Normalize salary within each department
df['Normalized Salary'] = df.groupby('Department')['Salary'].transform(lambda x:
(x - x.mean()) / x.std())
print(df)
```

# 🔗 2. Merging, Joining, Concatenating, and Reshaping Data

Combining datasets is essential for data wrangling. Pandas offers flexible functions for merging and joining datasets.

## ⊠ 2.1 Merging DataFrames with `merge()`

```python
# Sample DataFrames
df_employees = pd.DataFrame({
    'Employee': ['Alice', 'Bob', 'Charlie', 'David'],
    'Department': ['HR', 'HR', 'IT', 'IT']
})

df_salaries = pd.DataFrame({
    'Employee': ['Alice', 'Bob', 'Charlie', 'David'],
    'Salary': [70000, 80000, 90000, 85000]
})

# Merge on 'Employee'
merged_df = pd.merge(df_employees, df_salaries, on='Employee')
print(merged_df)
```

**Output:**

```
   Employee Department  Salary
0    Alice          HR   70000
1      Bob          HR   80000
2  Charlie          IT   90000
3    David          IT   85000
```

## ⚖️ 2.2 Types of Joins

- **Inner Join** → Keeps only matching rows.
- **Left Join** → Keeps all rows from the left DataFrame.
- **Right Join** → Keeps all rows from the right DataFrame.
- **Outer Join** → Keeps all rows from both DataFrames.

```python
# Outer join example
df1 = pd.DataFrame({'ID': [1, 2, 3], 'Name': ['Alice', 'Bob', 'Charlie']})
df2 = pd.DataFrame({'ID': [2, 3, 4], 'Salary': [80000, 90000, 85000]})

outer_join = pd.merge(df1, df2, on='ID', how='outer')
print(outer_join)
```

**Output:**

```
    ID     Name    Salary
0   1    Alice        NaN
1   2      Bob    80000.0
2   3  Charlie    90000.0
3   4      NaN    85000.0
```

---

## 📦 2.3 Concatenating DataFrames with `concat()`

```python
# Concatenate vertically
df1 = pd.DataFrame({'Name': ['Alice', 'Bob'], 'Age': [25, 30]})
df2 = pd.DataFrame({'Name': ['Charlie', 'David'], 'Age': [35, 40]})

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

**Output:**

```
      Name  Age
0    Alice   25
1      Bob   30
2  Charlie   35
3    David   40
```

---

## 🔁 2.4 Reshaping Data with `melt()` and `pivot()`

- `melt()` → Unpivot DataFrame from wide to long format.
- `pivot()` → Convert long format back to wide.

```python
# Melt example
df = pd.DataFrame({'Name': ['Alice', 'Bob'],
                   'Math': [85, 90],
                   'English': [78, 82]})

df_melted = pd.melt(df, id_vars=['Name'], var_name='Subject', value_name='Score')
print(df_melted)
```

**Output:**

```
    Name  Subject  Score
0  Alice     Math     85
1    Bob     Math     90
2  Alice  English     78
3    Bob  English     82
```

# 📊 3. Pivot Tables, Cross-Tabulations, and Complex Reshaping

## 🔢 3.1 Pivot Tables for Data Summarization

```python
df = pd.DataFrame({
    'Department': ['HR', 'HR', 'IT', 'IT', 'Finance', 'Finance'],
    'Gender': ['Female', 'Male', 'Male', 'Female', 'Female', 'Male'],
    'Salary': [70000, 80000, 90000, 85000, 75000, 72000]
})

# Create Pivot Table
pivot_table = pd.pivot_table(df, values='Salary', index='Department',
columns='Gender', aggfunc='mean')
print(pivot_table)
```

**Output:**

```
Gender       Female      Male
Department
Finance     75000.0   72000.0
HR          70000.0   80000.0
IT          85000.0   90000.0
```

## 📋 3.2 Cross-Tabulations with `crosstab()`

Useful for frequency tables.

```python
# Crosstab example
crosstab = pd.crosstab(df['Department'], df['Gender'])
print(crosstab)
```

**Output:**

```
Gender       Female  Male
Department
Finance           1     1
```

```
HR                1    1
IT                1    1
```

---

## 📏 3.3 MultiIndex Reshaping with `stack()` and `unstack()`

```python
# Create MultiIndex DataFrame
df_multi = df.set_index(['Department', 'Gender'])

# Stack the DataFrame
stacked = df_multi.stack()

# Unstack to reshape
unstacked = df_multi.unstack()
```

---

# 🗂 4. Handling Nested and Hierarchical Data

### 📦 4.1 Working with Nested JSON Data

```python
import json

# Sample nested JSON
data = '''
[
    {"name": "Alice", "details": {"age": 25, "city": "NY"}},
    {"name": "Bob", "details": {"age": 30, "city": "LA"}}
]
'''

# Load JSON
nested_json = json.loads(data)

# Normalize JSON
df_json = pd.json_normalize(nested_json, sep='_')
print(df_json)
```

**Output:**

```
    name  details_age details_city
0  Alice           25           NY
1    Bob           30           LA
```

---

### 🔖 4.2 Flattening Hierarchical Data

```python
# Multi-level column DataFrame
df = pd.DataFrame({
    ('Employee', 'Name'): ['Alice', 'Bob'],
    ('Employee', 'Age'): [25, 30],
    ('Salary', 'Base'): [70000, 80000],
    ('Salary', 'Bonus'): [5000, 6000]
})

# Flatten columns
df.columns = ['_'.join(col) for col in df.columns]
print(df)
```

**Output:**

```
  Employee_Name  Employee_Age  Salary_Base  Salary_Bonus
0         Alice            25        70000          5000
1           Bob            30        80000          6000
```

## 📌 Chapter Summary:

- Mastered **grouping**, **aggregating**, and applying **custom functions**.
- Efficiently **merged**, **joined**, and **concatenated** multiple datasets.
- Built complex **pivot tables**, **cross-tabs**, and used **reshaping tools**.
- Flattened and worked with **nested JSON** and **hierarchical data** structures.

# 📖 Chapter 9: Scaling Pandas for Big Data

Pandas is powerful for data manipulation but struggles with extremely large datasets that exceed system memory. In this chapter, we'll explore strategies to scale Pandas for **big data** workflows, including memory optimizations, out-of-core processing, and leveraging libraries like **Dask** and **Vaex** for parallel and distributed computing.

---

## 💾 1. Working with Large Datasets (Out-of-Core Processing)

Out-of-core processing allows you to work with datasets that don't fit into memory by loading data in **chunks**.

### 🗃️ 1.1 Reading Data in Chunks

Use the `chunksize` parameter in Pandas' `read_csv()` to process large files incrementally.

```python
import pandas as pd

# Load CSV in chunks of 100,000 rows
chunksize = 100_000
chunk_iter = pd.read_csv('large_dataset.csv', chunksize=chunksize)

# Process chunks
for chunk in chunk_iter:
    chunk_result = chunk.groupby('Category')['Sales'].sum()
    print(chunk_result.head())
```

This reads and processes the file in manageable pieces rather than loading it all at once.

---

### 🗄️ 1.2 Incremental Aggregation Example

If you want to compute global statistics across chunks:

```python
# Initialize empty DataFrame
final_result = pd.Series(dtype='float64')

# Aggregate across chunks
for chunk in pd.read_csv('large_dataset.csv', chunksize=chunksize):
    chunk_result = chunk.groupby('Region')['Sales'].sum()
    final_result = final_result.add(chunk_result, fill_value=0)

print(final_result)
```

---

# 🧠 2. Optimizing Memory Usage and Performance

Handling big data efficiently requires reducing memory consumption.

## 📏 2.1 Reducing Memory Footprint

### 2.1.1 Check Memory Usage

```python
df = pd.read_csv('large_dataset.csv')
print(df.info(memory_usage='deep'))
```

### 2.1.2 Downcasting Data Types

Convert larger data types to more efficient ones where possible.

```python
# Convert float64 to float32 and int64 to int32
df['Sales'] = pd.to_numeric(df['Sales'], downcast='float')
df['Quantity'] = pd.to_numeric(df['Quantity'], downcast='integer')
```

### 2.1.3 Convert Object Columns to Category

If a column has a limited number of unique values, convert it to `category` to save space.

```python
df['Product'] = df['Product'].astype('category')
```

💡 **Tip:** This is especially useful for columns like **Product IDs**, **Countries**, or **Categories**.

## 🚀 2.2 Efficient Data Loading

### 2.2.1 Selective Column Loading

Only load necessary columns using the `usecols` parameter.

```python
df = pd.read_csv('large_dataset.csv', usecols=['Date', 'Product', 'Sales'])
```

### 2.2.2 Parse Dates While Reading

Parsing dates during read avoids post-processing.

```python
df = pd.read_csv('large_dataset.csv', parse_dates=['OrderDate'])
```

# ⚡ 3. Using Dask and Vaex for Parallel and Distributed DataFrames

When Pandas is not enough, use **Dask** or **Vaex** for parallel processing and out-of-core data analysis.

## 💡 3.1 Dask: Parallel Computing with Pandas Syntax

Dask offers a parallel DataFrame similar to Pandas but optimized for large datasets.

### 3.1.1 Dask Setup and Basics

```python
import dask.dataframe as dd

# Load large CSV with Dask
dask_df = dd.read_csv('large_dataset.csv')

# Perform aggregation
result = dask_df.groupby('Region')['Sales'].mean().compute()
print(result)
```

- **Lazy Evaluation:** Operations are queued until `.compute()` is called.
- **Parallelization:** Dask splits the data into partitions and processes them in parallel.

### 3.1.2 Optimizing Dask Workflows

- Use `persist()` to cache data in memory across computations.
- Monitor performance using Dask's **dashboard** (`dask.distributed`).

```python
from dask.distributed import Client

client = Client()  # Starts a local Dask cluster
```

## ⚡ 3.2 Vaex: Fast, Out-of-Core DataFrames

Vaex handles **billions of rows** efficiently without loading data into memory.

### 3.2.1 Vaex Setup and Usage

```python
import vaex
```

```
# Load large CSV with Vaex
vaex_df = vaex.from_csv('large_dataset.csv', convert=True)

# Perform lazy evaluation
result = vaex_df.groupby('Region', agg=vaex.agg.mean('Sales'))
print(result)
```

- **Zero-Copy Memory Mapping:** Fast file access without loading the entire dataset.
- **Lazy Evaluation:** Only computes when needed.

---

# 🗄 4. Chunking and Batch Processing Large Data Files

For massive datasets stored in flat files or databases, processing data in **batches** ensures stability.

---

## 📂 4.1 Writing Large DataFrames in Chunks

When saving large data, write in chunks to avoid memory issues.

```
# Write large DataFrame in chunks
for i, chunk in enumerate(pd.read_csv('large_dataset.csv', chunksize=50000)):
    chunk.to_csv(f'chunk_{i}.csv', index=False)
```

---

## 📊 4.2 Batch Processing with SQL Databases

When data is stored in databases, read and process data in batches using **SQLAlchemy**.

```
from sqlalchemy import create_engine

# Connect to database
engine = create_engine('sqlite:///sales_data.db')

# Read in batches
query = 'SELECT * FROM sales_data'
chunks = pd.read_sql(query, engine, chunksize=100000)

for chunk in chunks:
    print(chunk['Sales'].sum())
```

---

## 📄 4.3 Handling Compressed and Remote Data

Pandas can directly read compressed and remote files.

```
# Reading compressed CSV
df = pd.read_csv('data.csv.gz', compression='gzip')

# Reading data from a URL
df = pd.read_csv('https://example.com/large_dataset.csv')
```

## 📌 Chapter Summary

☑ Processed large datasets using **out-of-core** techniques. ☑ Optimized memory usage with **downcasting** and **categorical data types**. ☑ Leveraged **Dask** and **Vaex** for parallel and distributed computations. ☑ Used **chunking** and **batch processing** for working with large files.

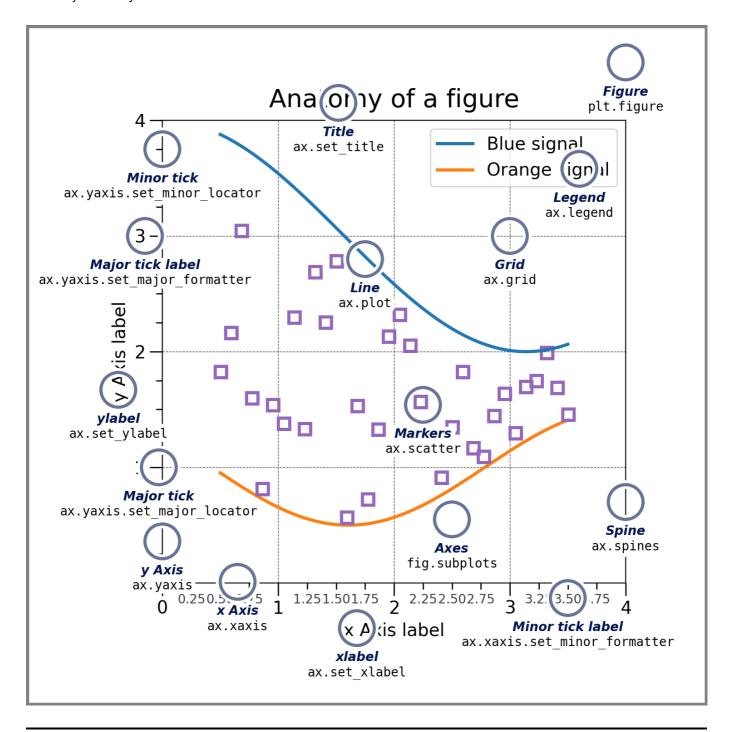# 📖 Chapter 10: Fundamentals of Data Visualization

Data visualization is a crucial part of data analysis, enabling you to interpret complex datasets through clear and impactful graphics. In this chapter, we'll dive into **Matplotlib**, the foundational plotting library in Python, covering its core concepts, essential plots, and powerful customization options.

## 🎨 1. Anatomy of a Matplotlib Plot

Before diving into creating visualizations, it's important to understand the structure of a Matplotlib plot.

### 🔍 1.1 Core Components

- **Figure** → The entire plotting window (container for one or more plots).
- **Axes** → The actual plot area inside the figure (can have multiple axes).
- **Axis** → The x-axis and y-axis of a plot.
- **Ticks** → Markers on the axes indicating specific data points.
- **Labels** → Text annotations for axes and data points.
- **Legend** → Describes data series within the plot.

Anatomy of a figure

## ▦ 1.2 Creating a Basic Plot

```python
import matplotlib.pyplot as plt
import numpy as np

# Sample Data
x = np.linspace(0, 10, 100)
y = np.sin(x)

# Create Figure and Axes
fig, ax = plt.subplots()

# Plotting
ax.plot(x, y, label='Sine Wave', color='blue')
```

```python
# Add Labels and Title
ax.set_title('Sine Wave Plot')
ax.set_xlabel('X-axis')
ax.set_ylabel('Y-axis')
ax.legend()

# Show Plot
plt.show()
```

# 📊 2. Creating Line Plots, Bar Charts, Histograms, and Scatter Plots

Let's dive into the most common types of visualizations used in data analysis.

## ☑ 2.1 Line Plots

Line plots are perfect for showing trends over time or continuous data.

```python
x = np.linspace(0, 10, 100)
y1 = np.sin(x)
y2 = np.cos(x)

plt.figure(figsize=(8,5))
plt.plot(x, y1, label='Sine', color='blue', linestyle='--')
plt.plot(x, y2, label='Cosine', color='red')
plt.title('Sine and Cosine Functions')
plt.xlabel('X-axis')
plt.ylabel('Y-axis')
plt.legend()
plt.show()
```

## 📊 2.2 Bar Charts

Bar charts are used for comparing categorical data.

```python
categories = ['A', 'B', 'C', 'D']
values = [23, 45, 56, 12]

plt.bar(categories, values, color='teal')
plt.title('Category Comparison')
plt.xlabel('Categories')
plt.ylabel('Values')
plt.show()
```

- Use `plt.barh()` for **horizontal** bar charts.

## ◻ 2.3 Histograms

Histograms visualize the distribution of numerical data.

```python
data = np.random.randn(1000)

plt.hist(data, bins=30, color='coral', edgecolor='black')
plt.title('Data Distribution')
plt.xlabel('Value')
plt.ylabel('Frequency')
plt.show()
```

- **Tip:** Adjust `bins` to control granularity.

## ◎ 2.4 Scatter Plots

Scatter plots are ideal for examining relationships between variables.

```python
x = np.random.rand(100)
y = np.random.rand(100)
sizes = 300 * np.random.rand(100)  # Varying bubble sizes
colors = np.random.rand(100)

plt.scatter(x, y, s=sizes, c=colors, alpha=0.6, cmap='viridis')
plt.title('Scatter Plot Example')
plt.xlabel('X-axis')
plt.ylabel('Y-axis')
plt.colorbar(label='Color Intensity')
plt.show()
```

- Add a **colorbar** for extra data dimensions.

# ⚒ 3. Customizing Visuals: Colors, Labels, Grids, and Annotations

Customizations can make your plots clearer and more engaging.

## ◔ 3.1 Colors and Styles

- Use color names (`'red'`), HEX codes (`'#FF5733'`), or color maps.
- Apply line styles: `'-'` (solid), `'--'` (dashed), `':'` (dotted).

```python
x = np.linspace(0, 10, 100)
y = np.sin(x)
```

```python
plt.plot(x, y, color='#1f77b4', linestyle='--', linewidth=2, marker='o')
plt.title('Customized Line Plot')
plt.show()
```

## 🏷 3.2 Adding Labels and Legends

- `set_xlabel()` and `set_ylabel()` for axes labels.
- `legend()` for data series.

```python
plt.plot(x, y, label='Sine Wave')
plt.xlabel('Time')
plt.ylabel('Amplitude')
plt.title('Labeled Plot')
plt.legend(loc='upper right')
plt.show()
```

## 📏 3.3 Grids and Annotations

- Add grids for better readability: `plt.grid(True)`
- Use annotations to highlight specific points.

```python
plt.plot(x, y)
plt.annotate('Peak', xy=(np.pi/2, 1), xytext=(2, 1.5),
             arrowprops=dict(facecolor='black', shrink=0.05))
plt.grid(True)
plt.show()
```

# 🖼 4. Subplots, Multiple Axes, and Plot Styling

For comparative analysis, you can create complex layouts with multiple plots.

## 📋 4.1 Creating Subplots

Use `plt.subplots()` to create multi-plot grids.

```python
fig, axs = plt.subplots(2, 2, figsize=(10, 8))

x = np.linspace(0, 10, 100)

axs[0, 0].plot(x, np.sin(x), 'r')
axs[0, 0].set_title('Sine')
```

```python
axs[0, 1].plot(x, np.cos(x), 'b')
axs[0, 1].set_title('Cosine')

axs[1, 0].plot(x, np.tan(x), 'g')
axs[1, 0].set_title('Tangent')

axs[1, 1].hist(np.random.randn(1000), bins=30, color='purple')
axs[1, 1].set_title('Histogram')

plt.tight_layout()
plt.show()
```

## 📊 4.2 Twin Axes for Different Scales

Overlay plots with different y-axis scales.

```python
fig, ax1 = plt.subplots()

x = np.arange(0, 10, 0.1)
y1 = np.sin(x)
y2 = y1**2

# Primary axis
ax1.plot(x, y1, 'b-')
ax1.set_xlabel('X-axis')
ax1.set_ylabel('Sine', color='b')

# Secondary axis
ax2 = ax1.twinx()
ax2.plot(x, y2, 'r--')
ax2.set_ylabel('Sine Squared', color='r')

plt.title('Twin Axes Example')
plt.show()
```

## 🎨 4.3 Applying Plot Styles

Matplotlib comes with pre-defined styles like ggplot, seaborn, and fivethirtyeight.

```python
plt.style.use('ggplot')

x = np.linspace(0, 10, 100)
y = np.sin(x)

plt.plot(x, y)
plt.title('Styled Plot with ggplot Theme')
plt.show()
```

View all styles:

```
print(plt.style.available)
```

---

## ☑️ **Chapter Summary**

◈ Mastered the **anatomy** of Matplotlib plots. ◈ Created **line plots, bar charts, histograms,** and **scatter plots**. ◈ Applied **customizations** for enhanced clarity and aesthetics. ◈ Built complex layouts with **subplots** and **multiple axes**.

---

# 📖 Chapter 11: Advanced Visualization Techniques

In this chapter, we'll explore advanced data visualization tools and techniques that go beyond basic plotting. From **statistical visualizations** using **Seaborn** to **interactive plots** with **Plotly**, and from **correlation heatmaps** to **geospatial** and **3D plots**, these methods will enable you to present complex data more intuitively and effectively.

## 🧮 1. Seaborn for Statistical Visualizations

**Seaborn** is built on top of Matplotlib, providing a high-level interface for drawing attractive and informative statistical graphics.

### 💡 1.1 Seaborn Setup

```python
import seaborn as sns
import matplotlib.pyplot as plt
import pandas as pd

# Load Seaborn's built-in dataset
tips = sns.load_dataset('tips')
```

## 📊 1.2 Enhanced Categorical Plots

### 1.2.1 Barplot with Confidence Intervals

```python
sns.barplot(x='day', y='total_bill', data=tips, ci='sd', palette='pastel')
plt.title('Average Total Bill per Day')
plt.show()
```

### 1.2.2 Boxplot and Violin Plot for Distribution Analysis

```python
fig, axs = plt.subplots(1, 2, figsize=(12, 5))

sns.boxplot(x='day', y='total_bill', data=tips, ax=axs[0], palette='Blues')
axs[0].set_title('Boxplot: Total Bill per Day')

sns.violinplot(x='day', y='total_bill', data=tips, ax=axs[1], palette='Oranges')
axs[1].set_title('Violin Plot: Total Bill per Day')

plt.tight_layout()
plt.show()
```

## 🌀 1.3 Pairplot for Exploratory Data Analysis

Quickly explore relationships between numerical variables:

```python
sns.pairplot(tips, hue='sex', palette='coolwarm')
plt.suptitle('Pairplot: Tips Dataset', y=1.02)
plt.show()
```

## ▒ 1.4 Heatmaps for Correlation Analysis

```python
corr = tips.corr(numeric_only=True)
sns.heatmap(corr, annot=True, cmap='YlGnBu', linewidths=0.5)
plt.title('Correlation Heatmap')
plt.show()
```

# 🖱 2. Plotly for Interactive and Web-Based Visualizations

**Plotly** creates highly interactive and web-friendly visualizations with minimal code.

## ⚡ 2.1 Installing Plotly

```
pip install plotly
```

## 📊 2.2 Interactive Line and Bar Plots

```python
import plotly.express as px

# Sample Data
df = px.data.gapminder().query("year == 2007")

# Interactive Scatter Plot
fig = px.scatter(df, x="gdpPercap", y="lifeExp", size="pop", color="continent",
                 hover_name="country", log_x=True, size_max=60, title="GDP vs Life
Expectancy")
fig.show()
```

Hover, zoom, and pan with ease in Plotly plots.

## 🌐 2.3 Choropleth Maps (Geospatial Visualizations)

```python
fig = px.choropleth(df, locations="iso_alpha",
                        color="lifeExp",
                        hover_name="country",
                        color_continuous_scale=px.colors.sequential.Plasma)
fig.update_layout(title="Global Life Expectancy (2007)")
fig.show()
```

## 📊 2.4 3D Scatter Plots

```python
fig = px.scatter_3d(df, x='gdpPercap', y='lifeExp', z='pop',
                        color='continent', size='pop', hover_name='country')
fig.update_layout(title='3D GDP vs Life Expectancy vs Population')
fig.show()
```

# 🔥 3. Heatmaps, Correlation Matrices, and Geospatial Plots

Advanced analysis often requires specialized visualizations.

## 🎛 3.1 Complex Correlation Heatmaps

Using Seaborn with mask options for cleaner heatmaps:

```python
import numpy as np

# Compute correlation
corr = tips.corr(numeric_only=True)

# Generate a mask for the upper triangle
mask = np.triu(np.ones_like(corr, dtype=bool))

# Create heatmap
sns.heatmap(corr, mask=mask, annot=True, cmap='coolwarm', linewidths=0.5)
plt.title('Upper Triangle Correlation Heatmap')
plt.show()
```

## 🗺 3.2 Geospatial Plots with Plotly and Geopandas

**Geopandas** allows manipulation of geospatial data, and Plotly enhances visualization.

```python
import geopandas as gpd
import plotly.express as px

# Load world map data
world = gpd.read_file(gpd.datasets.get_path('naturalearth_lowres'))

# Plot interactive map
fig = px.choropleth(world, locations='iso_a3', color='pop_est',
                    hover_name='name', color_continuous_scale='Viridis')
fig.update_layout(title='World Population by Country')
fig.show()
```

# 👦💻 4. 3D Plots, Network Graphs, and Complex Representations

## 🪀 4.1 Matplotlib 3D Plots

```python
from mpl_toolkits.mplot3d import Axes3D

fig = plt.figure(figsize=(8,6))
ax = fig.add_subplot(111, projection='3d')

# Data
x = np.random.rand(100)
y = np.random.rand(100)
z = np.random.rand(100)

# 3D Scatter
ax.scatter(x, y, z, c=z, cmap='viridis')
ax.set_title('3D Scatter Plot')
plt.show()
```

## 🕸️ 4.2 Network Graphs with NetworkX

**NetworkX** enables graph-based data visualizations like social networks and relationships.

```python
import networkx as nx

# Create graph
G = nx.erdos_renyi_graph(n=20, p=0.2)

# Plot graph
plt.figure(figsize=(10, 7))
nx.draw(G, with_labels=True, node_color='skyblue', node_size=700,
edge_color='gray')
plt.title('Random Network Graph')
plt.show()
```

## ⌖ 4.3 Sankey Diagrams for Flow Visualizations (Plotly)

```python
import plotly.graph_objects as go

fig = go.Figure(data=[go.Sankey(
    node=dict(label=["Source A", "Source B", "Target C", "Target D"]),
    link=dict(source=[0, 1, 0, 1],
              target=[2, 2, 3, 3],
              value=[8, 4, 2, 8]))])

fig.update_layout(title_text="Sankey Diagram Example", font_size=12)
fig.show()
```

# ✅ Chapter Summary

◈ Enhanced statistical visualizations with **Seaborn** ◈ Built **interactive** and **web-based** visualizations using **Plotly** ◈ Explored **heatmaps**, **correlation matrices**, and **geospatial plots** ◈ Created **3D plots** and **network graphs** for complex data analysis

# 📖 Chapter 12: Building Interactive Dashboards

In this chapter, we will explore how to create **interactive dashboards** that combine the power of **Pandas**, **Matplotlib**, and **Plotly**. We'll also integrate **ipywidgets** and **Voila** to add interactivity, and conclude with methods for creating reproducible reports using **Jupyter** and **nbconvert**.

By the end of this chapter, you'll be able to:

☑ Build dashboards that update dynamically based on user input ☑ Design visually engaging, interactive data applications ☑ Convert Jupyter Notebooks into standalone dashboards

## 🏗 1. Integrating Pandas, Matplotlib, and Plotly for Dashboards

To build a powerful dashboard, we'll combine data processing (**Pandas**), static visualizations (**Matplotlib**), and interactive plots (**Plotly**).

### 📊 1.1 Building a Simple Dashboard with Pandas & Matplotlib

**Step 1: Prepare the Data**

```python
import pandas as pd
import matplotlib.pyplot as plt

# Sample dataset
df = pd.read_csv('https://raw.githubusercontent.com/mwaskom/seaborn-data/master/iris.csv')

# Quick data check
df.head()
```

**Step 2: Create a Dynamic Plot**

```python
def plot_species_distribution(species):
    filtered_df = df[df['species'] == species]

    plt.figure(figsize=(8, 5))
    plt.hist(filtered_df['sepal_length'], bins=10, color='skyblue', edgecolor='black')
    plt.title(f'Sepal Length Distribution for {species.capitalize()}')
    plt.xlabel('Sepal Length')
    plt.ylabel('Frequency')
    plt.show()
```

```
# Example call
plot_species_distribution('setosa')
```

## ⚡ 1.2 Adding Interactivity with Plotly

```python
import plotly.express as px

# Interactive scatter plot
fig = px.scatter(df, x='sepal_width', y='sepal_length', color='species',
                 size='petal_length', hover_data=['petal_width'])
fig.update_layout(title='Iris Dataset: Sepal Dimensions')
fig.show()
```

Plotly enables zooming, panning, and tooltips without extra code.

# ▦ 2. Using ipywidgets and Voila for Interactivity

**ipywidgets** lets you create interactive controls, while **Voila** turns Jupyter notebooks into web apps.

### 🗎 2.1 Install ipywidgets and Voila

```
pip install ipywidgets voila
```

### ▦ 2.2 Creating Interactive Controls with ipywidgets

```python
import ipywidgets as widgets
from IPython.display import display

# Dropdown widget
species_dropdown = widgets.Dropdown(
    options=df['species'].unique(),
    value='setosa',
    description='Species:',
)

# Function to update plot
def update_plot(species):
    plot_species_distribution(species)

# Link widget to function
widgets.interactive(update_plot, species=species_dropdown)
```

Run this cell and select different species from the dropdown to see live plot updates.

---

## 🌐 2.3 Deploying Dashboards with Voila

Voila removes code cells and only displays widgets and outputs.

**Step 1: Save Your Notebook (e.g., `iris_dashboard.ipynb`)**

**Step 2: Run Voila**

```
voila iris_dashboard.ipynb
```

This launches a browser-based dashboard that anyone can interact with.

---

# 📄 3. Designing Reproducible Reports with Jupyter and nbconvert

To share your analysis as a report or presentation:

## 📦 3.1 Install nbconvert

```
pip install nbconvert
```

---

## 📄 3.2 Export as PDF or HTML

In Jupyter Notebook:

**File → Download As → HTML/PDF**

Or use terminal commands:

```
jupyter nbconvert --to html your_notebook.ipynb
jupyter nbconvert --to pdf your_notebook.ipynb
```

This creates polished, shareable reports directly from your analysis.

---

# 💡 4. Building a Real-World Interactive Dashboard

## 📊 4.1 Case Study: COVID-19 Data Dashboard

**Step 1: Import Data**

```python
covid_url = 'https://covid.ourworldindata.org/data/owid-covid-data.csv'
covid_df = pd.read_csv(covid_url, parse_dates=['date'])

# Filter recent data
covid_df = covid_df[covid_df['date'] >= '2023-01-01']
```

**Step 2: Create Interactive Country Selector**

```python
country_selector = widgets.Dropdown(
    options=covid_df['location'].unique(),
    value='United States',
    description='Country:',
)

def plot_covid_trend(country):
    data = covid_df[covid_df['location'] == country]
    fig = px.line(data, x='date', y='new_cases', title=f'COVID-19 Cases in
{country}')
    fig.show()

widgets.interactive(plot_covid_trend, country=country_selector)
```

**Step 3: Deploy with Voila**

```
voila covid_dashboard.ipynb
```

This creates a fully interactive COVID-19 dashboard that updates in real-time.

# ☑ Chapter Summary

◈ Built dynamic plots using **Pandas**, **Matplotlib**, and **Plotly** ◈ Created interactive widgets using **ipywidgets**
◈ Deployed dashboards as web apps using **Voila** ◈ Generated reproducible reports with **nbconvert**

# 📖 Chapter 13: Optimizing Data Pipelines

In data analysis, working with large datasets often leads to bottlenecks in **performance** and **memory usage**. This chapter focuses on optimizing data pipelines for speed and efficiency. We'll cover **profiling**, **vectorization**, **lazy evaluation**, handling **sparse data**, and **efficient disk I/O** strategies to build high-performance data pipelines.

## 🚀 1. Profiling and Benchmarking Data Operations

Before optimizing code, it's essential to identify bottlenecks using **profiling** and **benchmarking** tools.

### 📏 1.1 Profiling with `cProfile` and `pstats`

The `cProfile` module helps analyze code performance.

```python
import cProfile
import pandas as pd

def process_data():
    df = pd.DataFrame({'A': range(1000000), 'B': range(1000000)})
    df['C'] = df['A'] + df['B']
    return df

cProfile.run('process_data()')
```

For a more detailed report:

```
python -m cProfile -s time your_script.py
```

### ⏱️ 1.2 Benchmarking with `%timeit` in Jupyter

The `%timeit` magic command measures execution time with high precision.

```python
import numpy as np

arr = np.random.rand(1_000_000)
%timeit np.sqrt(arr)  # Fast, vectorized operation
```

### 📊 1.3 Visualizing Bottlenecks with `line_profiler`

Install and use `line_profiler` to identify slow lines of code:

```
pip install line_profiler
```

```
%load_ext line_profiler

def slow_function():
    total = 0
    for i in range(1000000):
        total += i
    return total

%lprun -f slow_function slow_function()
```

# ⚡ 2. Performance Tuning with Vectorization and Lazy Evaluation

## 💡 2.1 Vectorization: Replacing Loops with Numpy Operations

Vectorized operations in **Numpy** and **Pandas** avoid Python-level loops, dramatically improving performance.

### ✗ Inefficient (Using Loops):

```
arr = np.random.rand(1_000_000)
result = []

for val in arr:
    result.append(val ** 2)
```

### ☑ Efficient (Vectorized):

```
result = arr ** 2  # ~100x faster
```

## 🔀 2.2 Lazy Evaluation with Dask

**Dask** enables lazy evaluation for larger-than-memory datasets.

```
pip install dask
```

```
import dask.dataframe as dd

# Load a large CSV with Dask
df = dd.read_csv('large_dataset.csv')

# Lazy computation (not yet executed)
result = df.groupby('category')['value'].mean()

# Trigger computation
result.compute()
```

**Key Benefits:** ☑ Handles out-of-core datasets ☑ Parallelizes operations across cores

---

## ▦ 2.3 Applying Vectorization in Pandas

Use built-in Pandas functions instead of `apply` or loops:

✕ **Inefficient:**

```
df['new_col'] = df['value'].apply(lambda x: x * 2)
```

☑ **Efficient:**

```
df['new_col'] = df['value'] * 2
```

---

# ⬡ 3. Handling Sparse Data and Optimizing Data Types

## ▱ 3.1 Optimizing Data Types in Pandas

Reducing memory usage by optimizing column types.

```
df = pd.read_csv('large_dataset.csv')

# Check memory usage
print(df.info(memory_usage='deep'))

# Convert float64 to float32 and int64 to int32
df['int_column'] = df['int_column'].astype('int32')
df['float_column'] = df['float_column'].astype('float32')
```

**Tip:** Use `category` type for low-cardinality text columns.

```python
df['category_column'] = df['category_column'].astype('category')
```

## 🧮 3.2 Using Sparse Matrices for High-Dimensional Data

**Scipy's sparse matrices** help when most data points are zeros.

```python
from scipy import sparse
import numpy as np

dense_matrix = np.random.randint(0, 2, size=(1000, 1000))
sparse_matrix = sparse.csr_matrix(dense_matrix)

print(f"Memory usage of dense matrix: {dense_matrix.nbytes / 1e6} MB")
print(f"Memory usage of sparse matrix: {sparse_matrix.data.nbytes / 1e6} MB")
```

Sparse matrices significantly reduce memory usage in cases of sparse data.

# 💾 4. Memory Mapping and Efficient Disk I/O

## 🗄 4.1 Memory-Mapped Files with Numpy

**Memory mapping** enables access to data on disk as if it were in memory without loading it all at once.

```python
import numpy as np

# Create a memory-mapped file
data = np.memmap('data.dat', dtype='float32', mode='w+', shape=(10000, 10000))
data[:] = np.random.rand(10000, 10000)
del data  # Flush to disk

# Access a subset without loading the entire file
mapped_data = np.memmap('data.dat', dtype='float32', mode='r', shape=(10000,
10000))
print(mapped_data[5000, 5000])  # Access specific data point
```

## 📄 4.2 Efficient File Formats: Parquet and Feather

For large-scale data storage and retrieval, use **columnar storage formats** like **Parquet** and **Feather**.

```
pip install pyarrow fastparquet
```

```
# Save DataFrame as Parquet
df.to_parquet('data.parquet', engine='pyarrow')

# Read Parquet (faster than CSV)
df = pd.read_parquet('data.parquet', engine='pyarrow')
```

**Benefits:** ☑ Faster read/write times ☑ Built-in compression ☑ Optimized for column-wise operations

---

## 🗀 4.3 Chunking Large Files for Out-of-Core Processing

When dealing with massive CSV files:

```
chunksize = 10 ** 6  # Read 1 million rows at a time

for chunk in pd.read_csv('large_file.csv', chunksize=chunksize):
    # Perform aggregation on each chunk
    result = chunk.groupby('category')['value'].sum()
    print(result.head())
```

**Tip:** Combine with Dask for parallel processing.

---

# ☑ Chapter Summary

◈ Identified bottlenecks using **profiling** and **benchmarking** ◈ Improved performance with **vectorization** and **lazy evaluation** ◈ Reduced memory usage by **optimizing data types** and using **sparse matrices** ◈ Enhanced I/O efficiency with **memory mapping** and **columnar file formats**

---

# 📖 Chapter 14: Memory Management and Scalability

As datasets grow in size and complexity, efficient **memory management** and **scalability** become critical in data analysis. This chapter focuses on techniques for reducing memory footprints, working with **large-scale datasets** using tools like **Dask** and **PySpark**, implementing **parallel and distributed data analysis**, and handling **real-time analytics** using **caching** and **streaming data**.

By the end of this chapter, you will be able to:

☑ Optimize memory usage in **Pandas** and **Numpy** ☑ Work with massive datasets using **Dask** and **PySpark**
☑ Apply parallel and distributed processing techniques ☑ Implement real-time analytics with **caching** and **streaming**

---

## 🧮 1. Reducing Memory Footprint in Pandas and Numpy

Efficient memory management is key when working with large datasets in **Pandas** and **Numpy**.

---

### 📑 1.1 Memory Optimization in Pandas

#### 🔍 Step 1: Analyze Memory Usage

```python
import pandas as pd

# Load dataset
df = pd.read_csv('large_dataset.csv')

# Check memory usage
df.info(memory_usage='deep')
```

---

#### ⚡ Step 2: Optimize Data Types

- **Convert floats to float32** if precision allows.
- **Convert integers to smaller types** like int16 or int32.
- **Use categorical types** for repetitive text data.

```python
# Optimize numeric columns
df['int_column'] = pd.to_numeric(df['int_column'], downcast='integer')
df['float_column'] = pd.to_numeric(df['float_column'], downcast='float')

# Optimize string columns
df['category_column'] = df['category_column'].astype('category')
```

```python
# Check memory again
df.info(memory_usage='deep')
```

💡 **Tip:** These optimizations can reduce memory usage by **50-90%**.

---

## 📊 1.2 Efficient Memory Usage in Numpy

**Numpy** arrays can be optimized by choosing the right data types:

```python
import numpy as np

# Default: float64
arr = np.random.rand(1_000_000)

# Reduce memory with float32
arr_32 = arr.astype('float32')

print(f"Original size: {arr.nbytes / 1e6:.2f} MB")
print(f"Optimized size: {arr_32.nbytes / 1e6:.2f} MB")
```

**Result:** Float32 reduces memory by **50%** with minimal precision loss.

---

# 💻 2. Working with Large-Scale Datasets Using Dask and PySpark

When data exceeds your machine's memory, use distributed frameworks like **Dask** and **PySpark**.

---

## 🏗️ 2.1 Scaling Pandas with Dask

**Dask** extends Pandas for parallel and out-of-core computations.

📋 **Install Dask:**

```
pip install dask[complete]
```

🗒️ **Example: Using Dask DataFrames**

```python
import dask.dataframe as dd

# Load large CSV with Dask
df = dd.read_csv('large_dataset.csv')

# Lazy operations
result = df.groupby('category')['value'].mean()
```

```
# Execute computation
result.compute()
```

💡 **Key Benefits:** ☑ Handles datasets larger than RAM ☑ Parallelizes across multiple CPU cores ☑ Syntax similar to Pandas

---

## 🔥 2.2 Distributed Computing with PySpark

**PySpark** is a Python API for **Apache Spark**, designed for distributed data analysis across clusters.

📦 **Install PySpark:**

```
pip install pyspark
```

⚡ **Example: PySpark DataFrame**

```python
from pyspark.sql import SparkSession

# Initialize Spark
spark = SparkSession.builder.appName("BigDataApp").getOrCreate()

# Load data
df = spark.read.csv('large_dataset.csv', header=True, inferSchema=True)

# Group and aggregate
df.groupBy('category').avg('value').show()
```

💡 **Why PySpark?** ☑ Handles terabyte-scale datasets ☑ Built for distributed environments ☑ Supports SQL-like queries, machine learning, and graph processing

---

# 🧵 3. Parallel and Distributed Data Analysis Techniques

---

## ⚡ 3.1 Parallel Processing in Pandas Using `joblib`

For parallelizing custom functions in Pandas:

```
pip install joblib
```

```python
from joblib import Parallel, delayed
import pandas as pd
```

```python
df = pd.DataFrame({'A': range(1000000)})

# Parallel apply
def process_row(x):
    return x ** 2

results = Parallel(n_jobs=-1)(delayed(process_row)(x) for x in df['A'])
```

💡 **Tip:** `n_jobs=-1` uses all available CPU cores.

## 📊 3.2 Distributed Data Processing with Dask and PySpark

| Framework | Best For | Strength |
|-----------|----------|----------|
| **Dask** | Medium datasets, Pythonic | Easy Pandas migration |
| **PySpark** | Big data, distributed jobs | Cluster-level parallelization |

# 🌐 4. Caching, Streaming Data, and Real-Time Analytics

For **real-time data analysis** and **low-latency applications**, leverage caching and streaming frameworks.

## 🗄 4.1 Caching Data for Speed

**Pandas** and **Dask** can cache intermediate results to avoid redundant computations.

📇 **In Dask:**

```python
import dask.cache

# Enable cache
cache = dask.cache.Cache(1e9)  # 1 GB cache
cache.register()

# Cached computation
df = dd.read_csv('large_dataset.csv').persist()
```

## 📡 4.2 Streaming Data with PySpark Structured Streaming

Handle real-time data streams using PySpark:

```python
from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("StreamingApp").getOrCreate()
```

```python
# Read streaming data (e.g., from Kafka or socket)
stream_df = spark.readStream.format("socket") \
    .option("host", "localhost") \
    .option("port", 9999) \
    .load()

# Apply real-time transformations
query = stream_df.writeStream.outputMode("append").format("console").start()
query.awaitTermination()
```

💡 **Use Cases:** ☑ IoT device streams ☑ Real-time dashboards ☑ Log file monitoring

---

### ⚡ 4.3 Real-Time Analytics with Kafka and Spark

For complex real-time systems, combine **Kafka** (for data streams) with **Spark** (for processing).

**Pipeline Example:**

1. Kafka streams data →
2. Spark processes data →
3. Dashboards update in real-time

---

# ☑ Chapter Summary

◈ Reduced memory usage with efficient data types and sparse matrices ◈ Scaled data pipelines using **Dask** and **PySpark** ◈ Applied parallel and distributed data analysis techniques ◈ Implemented real-time analytics using **caching** and **streaming**

# 📖 Chapter 15: Advanced Data Wrangling and Transformation

In large-scale data analysis, mastering **data wrangling** and **transformation** is essential to extract meaningful insights. This chapter focuses on advanced techniques for reshaping data, performing complex aggregations, working with time series, and handling high-dimensional data.

By the end of this chapter, you will be able to:

☑ Reshape complex datasets using **melt**, **pivot**, **stack**, and **unstack** ☑ Apply advanced **GroupBy** operations, **rolling windows**, and **window functions** ☑ Perform **time series analysis** and detect anomalies ☑ Manage **high-dimensional** and **multivariate** data

## 📊 1. Complex Reshaping Techniques (Melt, Pivot, Stack, Unstack)

Reshaping data is key to aligning it with specific analysis needs. Pandas provides versatile tools like **melt**, **pivot**, **stack**, and **unstack**.

### 🔄 1.1 Melting DataFrames

The `melt` function transforms wide-format data into long-format, making it suitable for analysis.

📋 **Example:**

```python
import pandas as pd

df = pd.DataFrame({
    'Product': ['A', 'B', 'C'],
    'Jan_Sales': [200, 150, 300],
    'Feb_Sales': [220, 180, 320]
})

# Melt the DataFrame
melted_df = pd.melt(df, id_vars=['Product'], var_name='Month', value_name='Sales')

print(melted_df)
```

**Output:**

```
   Product      Month  Sales
0        A  Jan_Sales    200
1        B  Jan_Sales    150
2        C  Jan_Sales    300
3        A  Feb_Sales    220
```

```
4         B  Feb_Sales     180
5         C  Feb_Sales     320
```

---

## 🔁 1.2 Pivoting DataFrames

The **pivot** method reverses the effect of **melt**, turning long-format data back into wide format.

```python
pivot_df = melted_df.pivot(index='Product', columns='Month', values='Sales')
print(pivot_df)
```

**Output:**

```
Month     Jan_Sales  Feb_Sales
Product
A              200        220
B              150        180
C              300        320
```

---

## 📦 1.3 Stacking and Unstacking

- **stack()** converts columns into rows (creates a hierarchical index).
- **unstack()** moves rows back to columns.

📋 **Example:**

```python
df = pd.DataFrame({
    'City': ['NY', 'LA', 'SF'],
    '2023_Sales': [500, 600, 550],
    '2024_Sales': [520, 610, 580]
}).set_index('City')

# Stack: Convert columns into rows
stacked = df.stack()
print(stacked)

# Unstack: Reverse operation
unstacked = stacked.unstack()
print(unstacked)
```

---

# 📇 2. Advanced GroupBy, Rolling Windows, and Window Functions

## 🧮 2.1 Complex GroupBy Operations

Use **groupby** with custom functions and multi-level aggregations.

### 📋 Example: Multiple Aggregations

```python
df = pd.DataFrame({
    'Department': ['HR', 'HR', 'IT', 'IT', 'Finance', 'Finance'],
    'Employee': ['Alice', 'Bob', 'Charlie', 'David', 'Eve', 'Frank'],
    'Salary': [50000, 55000, 70000, 72000, 80000, 82000],
    'Bonus': [5000, 6000, 7000, 7500, 8000, 8500]
})

# Apply multiple aggregations
result = df.groupby('Department').agg({
    'Salary': ['mean', 'max'],
    'Bonus': 'sum'
})

print(result)
```

**Output:**

```
            Salary            Bonus
            mean      max     sum
Department
Finance     81000.0  82000   16500
HR          52500.0  55000   11000
IT          71000.0  72000   14500
```

## 🧮 2.2 Rolling Windows and Expanding Windows

**Rolling windows** and **expanding windows** help analyze time-series trends.

### 📋 Example: Rolling Mean

```python
dates = pd.date_range('2024-01-01', periods=10)
df = pd.DataFrame({'Date': dates, 'Sales': [5, 6, 7, 8, 7, 6, 5, 6, 7, 8]})
df.set_index('Date', inplace=True)

# Apply a rolling window
df['Rolling_Avg'] = df['Sales'].rolling(window=3).mean()
print(df)
```

### 📊 Expanding Window:

```python
df['Expanding_Mean'] = df['Sales'].expanding().mean()
```

## 📑 2.3 Window Functions with groupby

Apply window functions like **rank**, **lead/lag**, and **cumsum**.

### 📋 Example: Cumulative Sum

```python
df['Cumulative_Sum'] = df['Sales'].cumsum()
print(df)
```

# ⏰ 3. Time Series Analysis and Anomaly Detection

## 📆 3.1 Time Series Manipulation in Pandas

Use resample, shift, and rolling for time-based aggregations.

### 📋 Example: Monthly Resampling

```python
date_rng = pd.date_range(start='2024-01-01', end='2024-12-31', freq='D')
df = pd.DataFrame(date_rng, columns=['date'])
df['sales'] = np.random.randint(100, 500, size=(len(date_rng)))
df.set_index('date', inplace=True)

# Resample to monthly sales
monthly_sales = df.resample('M').sum()
print(monthly_sales.head())
```

## ⚠️ 3.2 Anomaly Detection Using Z-Score

Detect outliers using the **Z-Score** method.

```python
from scipy import stats

df['z_score'] = stats.zscore(df['sales'])

# Flag anomalies
df['anomaly'] = df['z_score'].apply(lambda x: 'Anomaly' if abs(x) > 2 else
'Normal')
```

## ☑ 3.3 Seasonal Decomposition with `statsmodels`

```
pip install statsmodels
```

```python
from statsmodels.tsa.seasonal import seasonal_decompose

result = seasonal_decompose(df['sales'], model='additive', period=30)
result.plot();
```

This decomposes time series into **trend**, **seasonality**, and **residuals**.

---

# 🧬 4. Handling High-Dimensional and Multivariate Data

---

## 📊 4.1 Reducing Dimensionality with PCA

For high-dimensional data, **Principal Component Analysis (PCA)** simplifies datasets while retaining most variance.

```python
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler
import numpy as np

# Simulated data
data = np.random.rand(100, 10)

# Standardize
scaler = StandardScaler()
data_scaled = scaler.fit_transform(data)

# Apply PCA
pca = PCA(n_components=2)
principal_components = pca.fit_transform(data_scaled)

print(principal_components[:5])  # First 5 data points
```

## ❄ 4.2 Handling Multivariate Time Series

Use **multi-indexing** and **panel data analysis** for time series with multiple variables.

### 📋 Example:

```python
arrays = [
    ['Region_A', 'Region_A', 'Region_B', 'Region_B'],
    pd.date_range('2024-01-01', periods=2).tolist() * 2
]

index = pd.MultiIndex.from_tuples(list(zip(*arrays)), names=['Region', 'Date'])

df = pd.DataFrame({'Sales': [100, 120, 150, 160]}, index=index)

# Access specific region data
print(df.loc['Region_A'])
```

## ☑ Chapter Summary

◈ Reshaped complex datasets using **melt**, **pivot**, and **stack/unstack** ◈ Applied advanced **GroupBy** techniques, rolling windows, and window functions ◈ Conducted **time series analysis** and detected anomalies ◈ Managed **high-dimensional data** using PCA and multi-indexing

# 📖 Chapter 16: Dimensionality Reduction and Complex Data Analysis

High-dimensional datasets can be overwhelming and challenging to interpret. Dimensionality reduction simplifies such data without significant loss of information, making it easier to visualize, analyze, and use in predictive models. This chapter covers essential techniques like **PCA**, **SVD**, **t-SNE**, and **UMAP**, along with methods for **feature engineering** and handling **multicollinearity**.

By the end of this chapter, you will be able to:

☑ Apply **PCA** and **SVD** for dimensionality reduction ☑ Use **t-SNE** and **UMAP** for data visualization ☑ Perform effective **feature engineering** for predictive analytics ☑ Analyze **correlations** and manage **multicollinearity**

---

## 📊 1. Principal Component Analysis (PCA) and Singular Value Decomposition (SVD)

### 🧮 1.1 Principal Component Analysis (PCA)

PCA is a linear dimensionality reduction technique that transforms data into new axes (principal components) that capture the most variance.

### 📋 Example: Applying PCA

```python
import pandas as pd
import numpy as np
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt

# Sample dataset
np.random.seed(42)
data = np.random.rand(100, 5) * 100
df = pd.DataFrame(data, columns=['Feature1', 'Feature2', 'Feature3', 'Feature4',
'Feature5'])

# Standardizing the data
scaler = StandardScaler()
scaled_data = scaler.fit_transform(df)

# Apply PCA
pca = PCA(n_components=2)
principal_components = pca.fit_transform(scaled_data)

# Create DataFrame for visualization
pca_df = pd.DataFrame(data=principal_components, columns=['PC1', 'PC2'])
```

```python
# Plot the results
plt.figure(figsize=(8, 6))
plt.scatter(pca_df['PC1'], pca_df['PC2'], alpha=0.7)
plt.title('PCA - First Two Principal Components')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.show()
```

⚡ **Key Points:**

- PCA reduces noise and redundancy.
- The explained variance ratio helps determine how much information is retained.

```python
print("Explained Variance Ratio:", pca.explained_variance_ratio_)
```

---

## ▦ 1.2 Singular Value Decomposition (SVD)

SVD is a matrix factorization method that breaks down a matrix into singular vectors and values. It's commonly used in PCA and recommendation systems.

📋 **Example: Using SVD**

```python
from numpy.linalg import svd

# Random matrix
A = np.random.rand(5, 5)

# Perform SVD
U, Sigma, Vt = svd(A)

print("U Matrix:\n", U)
print("\nSigma Values:\n", Sigma)
print("\nV Transposed:\n", Vt)
```

💡 **When to Use PCA vs. SVD:**

- **PCA** is optimized for variance preservation.
- **SVD** is versatile for dimensionality reduction, text mining (LSA), and matrix compression.

---

# 🗺 2. t-SNE, UMAP, and Other Visualization Techniques

For non-linear data structures, **t-SNE** and **UMAP** provide better clustering and visualization.

---

## 🌀 2.1 t-SNE (t-Distributed Stochastic Neighbor Embedding)

t-SNE is a non-linear technique for high-dimensional data visualization in 2D or 3D.

### 📋 Example: Visualizing Clusters with t-SNE

```python
from sklearn.manifold import TSNE

# Apply t-SNE
tsne = TSNE(n_components=2, perplexity=30, random_state=42)
tsne_results = tsne.fit_transform(scaled_data)

# Plot
plt.figure(figsize=(8, 6))
plt.scatter(tsne_results[:, 0], tsne_results[:, 1], c='blue', alpha=0.7)
plt.title('t-SNE Visualization')
plt.xlabel('Dimension 1')
plt.ylabel('Dimension 2')
plt.show()
```

### ⚡ Pros & Cons:

- ☑ Great for visualizing clusters.
- ✖ Computationally expensive on large datasets.

---

## 🌀 2.2 UMAP (Uniform Manifold Approximation and Projection)

UMAP preserves more global structure than t-SNE and is faster for large datasets.

### 📋 Example: Using UMAP

```python
import umap

# Apply UMAP
umap_reducer = umap.UMAP(n_neighbors=15, min_dist=0.1, random_state=42)
umap_results = umap_reducer.fit_transform(scaled_data)

# Plot
plt.figure(figsize=(8, 6))
plt.scatter(umap_results[:, 0], umap_results[:, 1], c='green', alpha=0.7)
plt.title('UMAP Visualization')
plt.xlabel('UMAP1')
plt.ylabel('UMAP2')
plt.show()
```

### ⚡ t-SNE vs. UMAP:

- **t-SNE** focuses on local relationships.
- **UMAP** preserves both local and global data structures.

---

# 🛠️ 3. Feature Engineering for Predictive Analytics

Feature engineering transforms raw data into meaningful features that improve model performance.

---

## ▦ 3.1 Creating New Features

- **Polynomial Features**: Capture non-linear relationships.
- **Aggregations**: Mean, sum, counts for groupings.
- **Domain Knowledge**: Derive features based on specific problem contexts.

📋 **Example: Polynomial Features**

```python
from sklearn.preprocessing import PolynomialFeatures

# Create interaction terms and polynomial features
poly = PolynomialFeatures(degree=2, include_bias=False)
poly_features = poly.fit_transform(df[['Feature1', 'Feature2']])

print("Polynomial Features Shape:", poly_features.shape)
```

---

## ⚖️ 3.2 Scaling and Normalization

- **StandardScaler** for Gaussian-like data.
- **MinMaxScaler** for bounded ranges.
- **RobustScaler** for outlier-heavy data.

```python
from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()
normalized_data = scaler.fit_transform(df)
```

---

## 🔄 3.3 Encoding Categorical Variables

- **One-Hot Encoding** for non-ordinal categories.
- **Ordinal Encoding** for ordered variables.

```python
df = pd.DataFrame({'Category': ['A', 'B', 'C', 'A', 'B']})
encoded_df = pd.get_dummies(df, columns=['Category'])
print(encoded_df)
```

# 📊 4. Correlation Analysis and Multicollinearity

Understanding relationships between variables is crucial in data analysis and modeling.

## 📏 4.1 Correlation Analysis

Use correlation matrices and heatmaps to visualize variable relationships.

### 📋 Example: Correlation Heatmap

```python
import seaborn as sns

# Correlation matrix
corr_matrix = df.corr()

# Plot heatmap
plt.figure(figsize=(8, 6))
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', linewidths=0.5)
plt.title('Feature Correlation Heatmap')
plt.show()
```

## ⚠️ 4.2 Handling Multicollinearity

Multicollinearity (high correlation between independent variables) can distort regression models.

### 💡 Detection Methods:

- **Variance Inflation Factor (VIF)**: Detects multicollinearity.

```python
from statsmodels.stats.outliers_influence import variance_inflation_factor

# Compute VIF for each feature
vif_data = pd.DataFrame()
vif_data["feature"] = df.columns
vif_data["VIF"] = [variance_inflation_factor(scaled_data, i) for i in range(scaled_data.shape[1])]

print(vif_data)
```

### ☑️ Solution Strategies:

- Remove highly correlated features.
- Apply dimensionality reduction (PCA).
- Use regularization techniques (Ridge, Lasso).

# ☑ Chapter Summary

◈ Applied **PCA** and **SVD** for linear dimensionality reduction ◈ Used **t-SNE** and **UMAP** for non-linear data visualization ◈ Performed **feature engineering** for predictive modeling ◈ Analyzed **correlations** and addressed **multicollinearity**

◈ Applied **PCA** and **SVD** for linear dimensionality reduction ◈ Used **t-SNE** and **UMAP** for non-linear data visualization ◈ Performed **feature engineering** for predictive modeling ◈ Analyzed **correlations** and addressed **multicollinearity**

# 📖 Chapter 17: Exploratory Data Analysis (EDA) in Practice

**Exploratory Data Analysis (EDA)** is a critical step in the data analysis pipeline. It helps analysts uncover patterns, detect anomalies, test hypotheses, and check assumptions using statistical summaries and visualizations. In this chapter, we will build effective EDA workflows, automate the process using custom Python scripts, and apply EDA techniques to a real-world case study.

By the end of this chapter, you will be able to:

☑ Design efficient **EDA workflows** ☑ Identify **patterns, trends, and anomalies** ☑ **Automate EDA** using Pythonic techniques ☑ Apply EDA in a **real-world case study**

## ◉ 1. Designing Effective EDA Workflows

A structured EDA process helps transform raw data into actionable insights. Here's a typical EDA workflow:

1. **Data Collection & Loading**
2. **Data Profiling & Quality Checks**
3. **Data Cleaning & Preprocessing**
4. **Exploratory Visualizations**
5. **Feature Engineering & Transformation**
6. **Pattern & Anomaly Detection**
7. **Insights Summary**

### ⚡ 1.1 Essential Python Libraries for EDA

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from pandas_profiling import ProfileReport
```

### 📋 1.2 Quick Data Overview

```python
# Sample E-commerce Dataset
df = pd.read_csv('ecommerce_data.csv')

# Quick data inspection
print(df.head())
print(df.info())
print(df.describe())
```

- `df.info()` → Data types, null values, memory usage
- `df.describe()` → Statistical summary of numerical features

💡 **Tip:** Use `pandas_profiling` for quick automated profiling.

```
profile = ProfileReport(df, title="E-Commerce Data Report", explorative=True)
profile.to_file("ecommerce_eda_report.html")
```

---

# 📊 2. Identifying Patterns, Trends, and Anomalies

Exploratory Visualizations help reveal hidden insights.

## ☑ 2.1 Univariate Analysis

Examine individual features.

```
# Distribution of Purchase Amount
sns.histplot(df['PurchaseAmount'], bins=30, kde=True)
plt.title('Distribution of Purchase Amount')
plt.show()
```

- **Histograms** → Distribution of numerical data
- **Bar plots** → Categorical data frequency

---

## 📊 2.2 Bivariate Analysis

Understand relationships between variables.

```
# Purchase Amount vs. Customer Age
sns.scatterplot(x='Age', y='PurchaseAmount', data=df)
plt.title('Age vs. Purchase Amount')
plt.show()

# Correlation Heatmap
plt.figure(figsize=(10, 6))
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
plt.title('Feature Correlation Matrix')
plt.show()
```

- **Scatterplots** → Continuous relationships
- **Heatmaps** → Correlations and multicollinearity

---

## ⚠ 2.3 Anomaly Detection

Outliers can distort analysis and predictions.

```python
# Boxplot to detect outliers
sns.boxplot(x=df['PurchaseAmount'])
plt.title('Outlier Detection in Purchase Amount')
plt.show()
```

**Advanced Techniques for Anomaly Detection:**

- **Z-Score Method**
- **IQR Method**
- **Isolation Forest & DBSCAN** (for complex data)

```python
from sklearn.ensemble import IsolationForest

# Apply Isolation Forest
iso_forest = IsolationForest(contamination=0.05)
df['anomaly'] = iso_forest.fit_predict(df[['PurchaseAmount', 'Age']])
sns.scatterplot(x='Age', y='PurchaseAmount', hue='anomaly', data=df)
plt.title('Anomaly Detection with Isolation Forest')
plt.show()
```

# 🤘 3. Automating EDA with Custom Scripts

For repetitive EDA tasks, Pythonic automation improves efficiency.

## 📋 3.1 Reusable EDA Function

```python
def automated_eda(dataframe):
    print("Shape of Dataset:", dataframe.shape)
    print("\nMissing Values:\n", dataframe.isnull().sum())
    print("\nData Types:\n", dataframe.dtypes)
    print("\nStatistical Summary:\n", dataframe.describe())

    # Correlation Heatmap
    plt.figure(figsize=(10, 6))
    sns.heatmap(dataframe.corr(), annot=True, cmap='Blues')
    plt.title('Correlation Matrix')
    plt.show()

# Run EDA
automated_eda(df)
```

## 🛠️ 3.2 Generating Automated Reports

```
# Full profiling report
profile = ProfileReport(df, title="Automated EDA Report", explorative=True)
profile.to_file("automated_eda_report.html")
```

## 📁 4. Case Study: Analyzing E-Commerce Customer Behavior

**Dataset:** Online retail transactions including demographics, purchase behavior, and browsing patterns.

🎯 **Objective:** Identify customer segments and factors affecting purchase amounts.

### 📋 4.1 Data Overview

```
# Load dataset
df = pd.read_csv('ecommerce_data.csv')

# Check for null values and data types
print(df.info())
```

### 📊 4.2 Key EDA Insights

1. **Who are the high-value customers?**
2. **Which products are most popular?**
3. **Are there seasonal trends in purchase patterns?**

### ☑ 4.3 Visualizing Customer Behavior

```
# Customer Segmentation by Age and Purchase Amount
sns.scatterplot(x='Age', y='PurchaseAmount', hue='CustomerSegment', data=df)
plt.title('Customer Segments: Age vs. Purchase Amount')
plt.show()

# Monthly Sales Trend
df['OrderDate'] = pd.to_datetime(df['OrderDate'])
df['Month'] = df['OrderDate'].dt.to_period('M')
monthly_sales = df.groupby('Month')['PurchaseAmount'].sum().reset_index()

plt.figure(figsize=(12, 6))
sns.lineplot(x='Month', y='PurchaseAmount', data=monthly_sales)
plt.title('Monthly Sales Trend')
plt.show()
```

### 📊 4.4 RFM (Recency, Frequency, Monetary) Analysis

**RFM Analysis** segments customers based on:

- **Recency** (How recently a customer made a purchase)
- **Frequency** (How often they purchase)
- **Monetary** (How much they spend)

```python
# RFM Calculation
snapshot_date = df['OrderDate'].max() + pd.Timedelta(days=1)
rfm = df.groupby('CustomerID').agg({
    'OrderDate': lambda x: (snapshot_date - x.max()).days,
    'InvoiceNo': 'count',
    'PurchaseAmount': 'sum'
})

rfm.rename(columns={'OrderDate': 'Recency',
                    'InvoiceNo': 'Frequency',
                    'PurchaseAmount': 'Monetary'}, inplace=True)

print(rfm.head())
```

**Visualization:**

```python
# RFM Distribution
fig, ax = plt.subplots(1, 3, figsize=(18, 5))
sns.histplot(rfm['Recency'], bins=30, ax=ax[0])
sns.histplot(rfm['Frequency'], bins=30, ax=ax[1])
sns.histplot(rfm['Monetary'], bins=30, ax=ax[2])
ax[0].set_title('Recency Distribution')
ax[1].set_title('Frequency Distribution')
ax[2].set_title('Monetary Distribution')
plt.show()
```

---

# ☑ Chapter Summary

◈ Designed structured **EDA workflows** ◈ Explored patterns, trends, and detected anomalies ◈ Automated EDA tasks using **Pythonic techniques** ◈ Applied EDA in a real-world **e-commerce case study**

---

# 📖 Chapter 18: Domain-Specific Data Analysis

In data analysis, domain knowledge shapes how we process, analyze, and interpret data. This chapter explores domain-specific data analysis techniques across finance, healthcare, social media, and scientific computing. We'll apply **Python libraries** like pandas, numpy, and visualization tools while integrating **domain-specific methodologies**.

By the end of this chapter, you'll be able to:

☑ Apply time series analysis for **financial data** ☑ Use statistical methods for **healthcare analytics** ☑ Perform text-based **sentiment analysis** ☑ Conduct **scientific simulations** using Monte Carlo methods

## 💰 1. Financial Data Analysis

Financial data often involves **time series**, where the order of data points is critical. Analysts focus on trends, seasonality, volatility, and risk.

### 📊 1.1 Time Series Analysis

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import yfinance as yf

# Load historical stock data
stock_data = yf.download('AAPL', start='2022-01-01', end='2023-12-31')
stock_data['Close'].plot(title='AAPL Stock Prices', figsize=(10, 5))
plt.show()
```

### ☑ 1.2 Moving Averages & Volatility

```python
# 20-day Moving Average
stock_data['20_MA'] = stock_data['Close'].rolling(window=20).mean()

# Bollinger Bands
stock_data['Upper'] = stock_data['20_MA'] + 2 *
stock_data['Close'].rolling(window=20).std()
stock_data['Lower'] = stock_data['20_MA'] - 2 *
stock_data['Close'].rolling(window=20).std()

# Visualization
plt.figure(figsize=(12, 6))
plt.plot(stock_data['Close'], label='Closing Price')
plt.plot(stock_data['20_MA'], label='20-Day MA', color='orange')
plt.fill_between(stock_data.index, stock_data['Upper'], stock_data['Lower'],
color='gray', alpha=0.3)
```

```python
plt.legend()
plt.title('AAPL Stock with Bollinger Bands')
plt.show()
```

## 📉 1.3 Forecasting with ARIMA

```python
from statsmodels.tsa.arima.model import ARIMA

# ARIMA Model
model = ARIMA(stock_data['Close'], order=(5, 1, 0))
model_fit = model.fit()
forecast = model_fit.forecast(steps=30)

# Plot Forecast
plt.figure(figsize=(10, 5))
plt.plot(stock_data['Close'], label='Actual')
plt.plot(pd.date_range(start=stock_data.index[-1], periods=30, freq='D'),
forecast, label='Forecast', color='red')
plt.legend()
plt.title('AAPL Stock Price Forecast')
plt.show()
```

## ⚡ 1.4 Risk Modeling: Value at Risk (VaR)

```python
# Historical VaR at 95% confidence
daily_returns = stock_data['Close'].pct_change().dropna()
var_95 = np.percentile(daily_returns, 5)
print(f"95% Value at Risk (VaR): {var_95 * 100:.2f}%")
```

---

# 🏥 2. Healthcare Analytics

Healthcare data requires handling sensitive, high-dimensional data with a focus on statistical accuracy.

## 📊 2.1 Patient Data Analysis

```python
# Sample healthcare data
df_health = pd.read_csv('patient_data.csv')

# Average patient age by disease type
df_health.groupby('Disease')['Age'].mean().plot(kind='bar', title='Average Age by
Disease')
plt.show()
```

## ☑ 2.2 Survival Analysis (Kaplan-Meier Curve)

```python
from lifelines import KaplanMeierFitter

# Sample survival data
df_survival = pd.DataFrame({
    'duration': [5, 6, 6, 2, 4, 4, 3],
    'event_observed': [1, 0, 0, 1, 1, 1, 0]
})

kmf = KaplanMeierFitter()
kmf.fit(df_survival['duration'], event_observed=df_survival['event_observed'])

kmf.plot()
plt.title('Kaplan-Meier Survival Curve')
plt.show()
```

## 🖩 2.3 Predictive Modeling for Disease Diagnosis

```python
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# Prepare data
X = df_health.drop('Diagnosis', axis=1)
y = df_health['Diagnosis']

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)

# Random Forest Classifier
model = RandomForestClassifier()
model.fit(X_train, y_train)
predictions = model.predict(X_test)

# Accuracy
print("Accuracy:", accuracy_score(y_test, predictions))
```

# 💬 3. Social Media & Text Data Analytics

Social media data is mostly **unstructured text**. Analyzing user sentiment and trends provides valuable insights.

## 📊 3.1 Text Preprocessing

```python
import pandas as pd
import re
from nltk.corpus import stopwords
```

```python
from nltk.tokenize import word_tokenize

# Sample social media data
tweets = pd.DataFrame({'text': ['I love AI! 😍', 'This is terrible... 😣',
'Excited for the future of tech. 🚀']})

# Text cleaning function
def clean_text(text):
    text = re.sub(r'[^A-Za-z0-9 ]+', '', text)  # Remove emojis & special chars
    tokens = word_tokenize(text.lower())
    return ' '.join([word for word in tokens if word not in
stopwords.words('english')])

tweets['cleaned'] = tweets['text'].apply(clean_text)
print(tweets)
```

## 😄 3.2 Sentiment Analysis Using NLTK

```python
from nltk.sentiment import SentimentIntensityAnalyzer

sia = SentimentIntensityAnalyzer()
tweets['sentiment_score'] = tweets['cleaned'].apply(lambda x:
sia.polarity_scores(x)['compound'])

# Classify sentiments
tweets['sentiment'] = tweets['sentiment_score'].apply(lambda x: 'Positive' if x >
0.05 else ('Negative' if x < -0.05 else 'Neutral'))
print(tweets[['text', 'sentiment']])
```

## ☑ 3.3 Word Cloud for Trend Analysis

```python
from wordcloud import WordCloud

# Combine all tweets
all_text = ' '.join(tweets['cleaned'])

# Generate word cloud
wordcloud = WordCloud(width=800, height=400,
background_color='white').generate(all_text)

plt.figure(figsize=(10, 5))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis('off')
plt.title('Social Media Word Cloud')
plt.show()
```

# 🎰 4. Scientific Computing: Monte Carlo Simulations

Monte Carlo simulations use **random sampling** to solve complex problems that are deterministic in nature but difficult to solve directly.

## 🎲 4.1 Estimating Pi with Monte Carlo

```python
import numpy as np
import matplotlib.pyplot as plt

# Simulation Parameters
n_points = 10000
circle_points = 0

# Generate random points
x = np.random.rand(n_points)
y = np.random.rand(n_points)

# Calculate distance from origin
distance = np.sqrt(x**2 + y**2)
circle_points = np.sum(distance <= 1)

# Estimate Pi
pi_estimate = (circle_points / n_points) * 4
print(f"Estimated Pi: {pi_estimate}")

# Visualization
plt.figure(figsize=(6, 6))
plt.scatter(x[distance <= 1], y[distance <= 1], color='blue', s=1, label='Inside Circle')
plt.scatter(x[distance > 1], y[distance > 1], color='red', s=1, label='Outside Circle')
plt.legend()
plt.title('Monte Carlo Estimation of Pi')
plt.show()
```

## 📊 4.2 Risk Analysis Using Monte Carlo

```python
# Simulate stock returns
simulations = []
for _ in range(1000):
    returns = np.random.normal(0.001, 0.02, 252)
    cumulative_return = np.prod(1 + returns) - 1
    simulations.append(cumulative_return)

# Plot results
plt.hist(simulations, bins=50, edgecolor='black')
plt.title('Monte Carlo Simulation of Annual Returns')
plt.xlabel('Return')
plt.ylabel('Frequency')
plt.show()
```

# ☑ Chapter Summary

◇ Explored **time series analysis** and **risk modeling** in finance ◇ Applied **statistical methods** and **predictive modeling** in healthcare ◇ Performed **text analysis** and **sentiment classification** for social media ◇ Used **Monte Carlo simulations** for complex scientific problems

# 📖 Chapter 19: Building Scalable Data Analysis Pipelines

Modern data projects require scalability, reproducibility, and maintainability. This chapter focuses on building **scalable data analysis pipelines** using **best practices** for structuring complex projects, implementing data versioning, automating ETL workflows, and deploying end-to-end pipelines.

By the end of this chapter, you'll be able to:

☑ Structure complex data projects for **scalability and reusability** ☑ Apply **data versioning** with DVC and Git ☑ Build **automated ETL (Extract, Transform, Load) pipelines** ☑ Implement a real-world **churn prediction pipeline**

---

## 📑 1. Structuring Complex Data Projects for Reusability

A well-structured project ensures maintainability, reproducibility, and collaboration.

### 📂 1.1 Recommended Project Structure

```
churn-prediction/
│
├── data/              # Raw and processed data
│   ├── raw/
│   └── processed/
│
├── notebooks/         # Jupyter notebooks for exploration
│
├── src/               # Source code for data pipelines
│   ├── data_ingestion.py
│   ├── data_preprocessing.py
│   ├── feature_engineering.py
│   └── model_training.py
│
├── models/            # Trained models
│
├── reports/           # Generated reports and visualizations
│
├── requirements.txt   # Dependencies
├── dvc.yaml           # Data pipeline configuration (DVC)
└── README.md
```

### 🛠 1.2 Best Practices

- **Modular Code:** Split code into reusable modules (`ingestion`, `preprocessing`, `modeling`).
- **Environment Management:** Use `requirements.txt` or `conda` for dependency control.
- **Reproducibility:** Leverage version control for code (`Git`) and data (`DVC`).

- **Automated Workflows:** Use tools like **Airflow** or **Prefect** for orchestrating tasks.

---

# 📊 2. Data Versioning and Reproducibility with DVC and Git

Traditional Git struggles with large datasets. **DVC (Data Version Control)** tracks datasets and models, ensuring reproducibility.

## ⚒️ 2.1 Setting Up DVC

1. **Install DVC:**

```
pip install dvc
```

2. **Initialize DVC in the Project:**

```
git init
dvc init
```

3. **Track Data Files:**

```
dvc add data/raw/customer_data.csv
git add data/.gitignore data/raw/customer_data.csv.dvc
git commit -m "Track raw customer data with DVC"
```

4. **Remote Storage (optional):**

```
dvc remote add -d myremote s3://mybucket/mydata
dvc push
```

## 🔄 2.2 Managing Data Versions

- **Modify Data:** Update your dataset (`data/raw/customer_data.csv`).
- **Track Changes:**

```
dvc add data/raw/customer_data.csv
git commit -am "Updated customer data"
dvc push
```

- **Switch Between Versions:**

```
git checkout previous_commit_hash
dvc checkout
```

💡 **Tip:** Combine **Git** (code) + **DVC** (data) for complete project reproducibility.

---

# 🔄 3. Creating Automated ETL Pipelines

An **ETL Pipeline** automates the process of **Extracting**, **Transforming**, and **Loading** data.

## ⚡ 3.1 Designing the ETL Flow

1. **Extract:** Pull data from CSV, APIs, or databases.
2. **Transform:** Clean, preprocess, and engineer features.
3. **Load:** Save processed data for modeling or analysis.

## 🛠️ 3.2 Building the ETL Pipeline in Python

```python
import pandas as pd
from sqlalchemy import create_engine

def extract_data(file_path):
    """Extract data from CSV."""
    return pd.read_csv(file_path)

def transform_data(df):
    """Data cleaning and feature engineering."""
    df.dropna(inplace=True)
    df['tenure_group'] = pd.cut(df['tenure'], bins=[0, 12, 24, 48, 72], labels=
['0-12', '12-24', '24-48', '48-72'])
    return df

def load_data(df, db_uri):
    """Load data into a SQL database."""
    engine = create_engine(db_uri)
    df.to_sql('processed_customers', engine, if_exists='replace', index=False)

# Run ETL
if __name__ == "__main__":
    raw_df = extract_data('data/raw/customer_data.csv')
    transformed_df = transform_data(raw_df)
    load_data(transformed_df, 'sqlite:///churn.db')
```

---

# 📊 4. Case Study: End-to-End Churn Prediction Pipeline

Let's build a complete churn prediction pipeline from data ingestion to model evaluation.

## 🗒️ 4.1 Step 1: Data Ingestion

```python
import pandas as pd

# Load customer data
df = pd.read_csv('data/raw/customer_data.csv')
print(df.head())
```

## 🔄 4.2 Step 2: Data Preprocessing

```python
from sklearn.preprocessing import LabelEncoder

# Encode categorical variables
le = LabelEncoder()
df['gender'] = le.fit_transform(df['gender'])

# Handle missing values
df.fillna(df.mean(), inplace=True)
```

## 📐 4.3 Step 3: Feature Engineering

```python
# Feature engineering example: Monthly spending ratio
df['spending_ratio'] = df['TotalCharges'] / (df['tenure'] + 1)
```

## 🤘 4.4 Step 4: Model Training

```python
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

# Train-test split
X = df.drop(['Churn'], axis=1)
y = df['Churn']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)

# Train model
model = RandomForestClassifier()
model.fit(X_train, y_train)

# Evaluate
predictions = model.predict(X_test)
accuracy = accuracy_score(y_test, predictions)
print(f"Model Accuracy: {accuracy:.2f}")
```

## ☑ 4.5 Step 5: Pipeline Automation with DVC

1. **Define Stages in `dvc.yaml`:**

```yaml
stages:
  preprocess:
    cmd: python src/data_preprocessing.py
    deps:
      - data/raw/customer_data.csv
      - src/data_preprocessing.py
    outs:
      - data/processed/customer_data.csv

  train:
    cmd: python src/model_training.py
    deps:
      - data/processed/customer_data.csv
      - src/model_training.py
    outs:
      - models/churn_model.pkl
```

2. **Run the Pipeline:**

```
dvc repro
```

## 📊 4.6 Step 6: Visualizing Churn Insights

```python
import seaborn as sns
import matplotlib.pyplot as plt

# Churn by Contract Type
sns.countplot(x='Contract', hue='Churn', data=df)
plt.title('Churn by Contract Type')
plt.show()
```

# ✅ Chapter Summary

◈ Structured data projects for **scalability** and **reusability** ◈ Applied **data versioning** with **DVC** ◈ Built **automated ETL pipelines** for efficient workflows ◈ Developed an **end-to-end churn prediction pipeline**

# 📖 Chapter 20: Integrating Data Analysis with External Systems

Modern data analysis often extends beyond local datasets, requiring integration with **databases**, **APIs**, **cloud platforms**, and **real-time streaming systems**. This chapter guides Python experts on how to bridge the gap between traditional data analysis and scalable, production-ready data workflows.

By the end of this chapter, you will be able to:

☑ Integrate data pipelines with **SQL/NoSQL databases** and **APIs** ☑ Perform data analysis using **cloud-based data warehouses** ☑ Process **real-time data streams** with Kafka and Pandas ☑ Apply **data engineering** techniques to optimize data workflows

## 🗄 1. Working with Databases (SQL, NoSQL) and APIs

Seamless integration with databases and external APIs is crucial for scalable data analysis.

### ▨ 1.1 Connecting to SQL Databases

**Libraries Used:** `SQLAlchemy`, `pandas`, `sqlite3`, `psycopg2`

#### 🔗 Example: Connecting to PostgreSQL

```python
import pandas as pd
from sqlalchemy import create_engine

# PostgreSQL connection string
engine =
create_engine('postgresql+psycopg2://user:password@localhost:5432/mydatabase')

# Query data into DataFrame
df = pd.read_sql("SELECT * FROM customers", con=engine)
print(df.head())
```

#### ⚡ Best Practices:

- Use **parameterized queries** to avoid SQL injection.
- Optimize queries for **large datasets** (use `LIMIT`, `WHERE` clauses).

### ▨ 1.2 Working with NoSQL Databases (MongoDB)

**Libraries Used:** `pymongo`, `pandas`

```python
from pymongo import MongoClient
import pandas as pd

# Connect to MongoDB
client = MongoClient("mongodb://localhost:27017/")
db = client["ecommerce"]
collection = db["orders"]

# Convert MongoDB data to DataFrame
data = list(collection.find())
df = pd.DataFrame(data)
print(df.head())
```

⚡ **Tips:**

- Use **indexes** for faster querying.
- Convert nested documents using `json_normalize()`.

---

🌐 **1.3 Consuming APIs for Data Extraction**

**Libraries Used:** `requests`, `pandas`

```python
import requests
import pandas as pd

# Sample API (e.g., weather data)
url = "https://api.weatherapi.com/v1/current.json"
params = {"key": "YOUR_API_KEY", "q": "London"}

response = requests.get(url, params=params)
data = response.json()

# Convert to DataFrame
df = pd.json_normalize(data)
print(df[['location.name', 'current.temp_c', 'current.condition.text']])
```

⚡ **Best Practices:**

- Handle **rate limits** and **timeouts**.
- Implement **error handling** for API failures.

---

☁ **2. Cloud Data Processing with AWS, BigQuery, and Snowflake**

Cloud data warehouses offer scalability, security, and performance for big data processing.

---

## ▨ 2.1 Using AWS S3 for Data Storage

**Libraries Used:** boto3, pandas

```python
import boto3
import pandas as pd
from io import StringIO

# Initialize S3 client
s3 = boto3.client('s3', aws_access_key_id='YOUR_KEY',
aws_secret_access_key='YOUR_SECRET')

# Download CSV from S3
response = s3.get_object(Bucket='mybucket', Key='data/sales.csv')
df = pd.read_csv(StringIO(response['Body'].read().decode('utf-8')))
print(df.head())
```

⚡ **Tips:**

- Use **S3 Select** for querying large files without downloading.
- Leverage **IAM roles** for secure access.

## ▨ 2.2 Analyzing Data in Google BigQuery

**Libraries Used:** google-cloud-bigquery, pandas

```python
from google.cloud import bigquery

# Initialize BigQuery client
client = bigquery.Client()

# SQL query
query = """
    SELECT customer_id, SUM(purchase_amount) as total_spent
    FROM `myproject.dataset.sales`
    GROUP BY customer_id
    ORDER BY total_spent DESC
    LIMIT 10
"""

# Convert query result to DataFrame
df = client.query(query).to_dataframe()
print(df)
```

⚡ **Best Practices:**

- Optimize queries with **partitioned** and **clustered** tables.

- Use **cost estimation** before running large queries.

---

## ▨ 2.3 Snowflake for Scalable Data Warehousing

**Libraries Used:** `snowflake-connector-python`, `pandas`

```python
import snowflake.connector
import pandas as pd

# Connect to Snowflake
conn = snowflake.connector.connect(
    user='username',
    password='password',
    account='account_name',
    warehouse='warehouse_name',
    database='database_name',
    schema='schema_name'
)

# Run query
cursor = conn.cursor()
cursor.execute("SELECT * FROM customers LIMIT 5")
df = pd.DataFrame(cursor.fetchall(), columns=[col[0] for col in
cursor.description])
print(df.head())
```

### ⚡ Tips:

- Use **external tables** for semi-structured data (JSON, Parquet).
- Scale compute power dynamically with **Snowflake Warehouses**.

---

# ♀ 3. Real-Time Data Streaming with Kafka and Pandas

Real-time data analysis is essential for industries like finance, IoT, and social media analytics.

---

## ▨ 3.1 Setting Up Kafka for Streaming

**Libraries Used:** `kafka-python`, `pandas`

1. **Install Kafka (local or cloud)**
2. **Start Producer and Consumer Services**

### ⚡ 3.2 Streaming Data from Kafka

```python
from kafka import KafkaConsumer
import json
```

```python
import pandas as pd

# Initialize Kafka consumer
consumer = KafkaConsumer(
    'realtime-data',
    bootstrap_servers='localhost:9092',
    auto_offset_reset='earliest',
    value_deserializer=lambda x: json.loads(x.decode('utf-8'))
)

# Stream data into DataFrame
data = []
for message in consumer:
    data.append(message.value)
    if len(data) >= 10:
        break

df = pd.DataFrame(data)
print(df.head())
```

⚡ **Best Practices:**

- Use **consumer groups** for load balancing.
- Implement **backpressure** handling for high-throughput streams.

---

# ⚙️ 4. Data Engineering for Data Science Workflows

Integrating **data engineering** practices ensures that data pipelines are efficient, maintainable, and production-ready.

---

## 🏛️ 4.1 Orchestrating Workflows with Apache Airflow

Airflow automates and manages complex workflows using Directed Acyclic Graphs (DAGs).

```python
from airflow import DAG
from airflow.operators.python import PythonOperator
from datetime import datetime

def extract_data():
    print("Extracting data...")

def transform_data():
    print("Transforming data...")

def load_data():
    print("Loading data...")

with DAG('etl_pipeline', start_date=datetime(2024, 1, 1),
    schedule_interval='@daily') as dag:
```

```
    extract = PythonOperator(task_id='extract', python_callable=extract_data)
    transform = PythonOperator(task_id='transform',
python_callable=transform_data)
    load = PythonOperator(task_id='load', python_callable=load_data)

    extract >> transform >> load
```

### ⚡ Best Practices:

- Use **XCom** for inter-task data sharing.
- Implement **failure recovery** and **retries**.

---

### ☑ 4.2 Building End-to-End ETL with Cloud Functions

**Example Workflow:**

1. **Trigger:** File upload to AWS S3
2. **Function:** AWS Lambda processes file
3. **Load:** Data is written to Redshift for analysis

---

# ✅ Chapter Summary

◈ Connected to **SQL/NoSQL databases** and **external APIs** ◈ Leveraged **cloud platforms** (AWS, BigQuery, Snowflake) for scalable analysis ◈ Streamed and processed **real-time data** using Kafka ◈ Implemented data engineering best practices with **Airflow** and **ETL pipelines**

---

# 📖 Chapter 21: From Data Analysis to Data Science

As data analysis evolves into more predictive and prescriptive forms, transitioning into **Data Science** becomes a natural progression. This chapter bridges the gap between data analysis and machine learning by introducing essential concepts and techniques for preparing data, building models, and evaluating their performance.

By the end of this chapter, you will be able to:

☑ Prepare data for machine learning models ☑ Use **Scikit-learn** for building predictive analytics pipelines ☑ Apply **feature engineering** with Pandas and Numpy ☑ Build and evaluate simple machine learning pipelines

## 🖌️ 1. Preparing Data for Machine Learning Models

The quality of your data directly impacts model performance. Preprocessing and feature engineering ensure that your data is suitable for machine learning algorithms.

### ⬤ 1.1 Data Preprocessing Workflow

Key preprocessing steps:

1. **Handling missing values**
2. **Encoding categorical variables**
3. **Scaling and normalizing features**
4. **Splitting data into training and test sets**

### 🔗 Example: Preparing a Dataset for ML

```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.impute import SimpleImputer
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline

# Sample dataset
df = pd.DataFrame({
    'age': [25, 30, 35, None, 40],
    'income': [50000, 60000, 70000, 80000, None],
    'city': ['NY', 'LA', 'SF', 'NY', 'LA'],
    'purchased': [0, 1, 0, 1, 0]
})

# Splitting features and target
X = df.drop('purchased', axis=1)
y = df['purchased']
```

```python
# Define preprocessors
numeric_features = ['age', 'income']
numeric_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='mean')),
    ('scaler', StandardScaler())
])

categorical_features = ['city']
categorical_transformer = Pipeline(steps=[
    ('encoder', OneHotEncoder())
])

# Combine preprocessors
preprocessor = ColumnTransformer(
    transformers=[
        ('num', numeric_transformer, numeric_features),
        ('cat', categorical_transformer, categorical_features)
    ])

# Apply preprocessing
X_preprocessed = preprocessor.fit_transform(X)
print(X_preprocessed)
```

## ⚡ 1.2 Best Practices:

- Always **scale numerical features** before feeding them into distance-based models (e.g., KNN, SVM).
- Use **OneHotEncoding** for nominal categories and **OrdinalEncoding** for ordered categories.
- Impute missing values carefully—consider domain knowledge.

# 🤘 2. Introduction to Scikit-learn for Predictive Analytics

**Scikit-learn** offers a wide range of machine learning algorithms and tools for building and evaluating predictive models.

## ▨ 2.1 Building a Simple Classification Model

Let's build a **Logistic Regression** model to predict customer purchases.

```python
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import Pipeline
from sklearn.metrics import accuracy_score

# Build pipeline with preprocessing and model
clf_pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('classifier', LogisticRegression())
])
```

```
# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Train model
clf_pipeline.fit(X_train, y_train)

# Predictions
y_pred = clf_pipeline.predict(X_test)

# Evaluate
accuracy = accuracy_score(y_test, y_pred)
print(f"Model Accuracy: {accuracy:.2f}")
```

◉ **Key Components:**

- **Pipeline:** Combines preprocessing and modeling in one streamlined process.
- **Train-Test Split:** Evaluates model performance on unseen data.
- **Accuracy Score:** Basic evaluation metric for classification.

---

## ▨ 2.2 Evaluating Model Performance

Besides accuracy, consider these metrics:

- **Precision & Recall:** Useful for imbalanced datasets
- **F1 Score:** Harmonic mean of precision and recall
- **ROC-AUC Curve:** Measures classifier's ability to distinguish between classes

```
from sklearn.metrics import classification_report, roc_auc_score

# Detailed report
print(classification_report(y_test, y_pred))

# ROC-AUC score
y_prob = clf_pipeline.predict_proba(X_test)[:, 1]
roc_auc = roc_auc_score(y_test, y_prob)
print(f"ROC-AUC Score: {roc_auc:.2f}")
```

---

# ▦ 3. Using Pandas and Numpy for Feature Engineering

**Feature engineering** is the art of extracting meaningful patterns from raw data.

---

## ◉ 3.1 Creating New Features

- **Date Features:** Extract day, month, or weekday from timestamps

- **Aggregations:** Group data and calculate summaries
- **Interactions:** Combine existing features to capture deeper relationships

🖇️ **Example: Feature Engineering with Pandas**

```python
import pandas as pd
import numpy as np

# Sample sales data
df = pd.DataFrame({
    'date': pd.date_range(start='2024-01-01', periods=5),
    'sales': [200, 300, 250, 400, 350],
    'store': ['A', 'B', 'A', 'B', 'A']
})

# Extracting date features
df['day_of_week'] = df['date'].dt.dayofweek

# Aggregating sales by store
store_sales = df.groupby('store')['sales'].sum().reset_index()

# Creating interaction term
df['sales_per_day'] = df['sales'] / (df['day_of_week'] + 1)

print(df)
```

## 🪩 3.2 Handling Outliers and Skewed Data

- **Log Transformations:** Reduce right skew in income or price data
- **Binning:** Group continuous data into discrete buckets
- **Winsorizing:** Limit extreme values to reduce the effect of outliers

```python
# Log transform
df['log_sales'] = np.log1p(df['sales'])

# Binning sales
df['sales_bin'] = pd.qcut(df['sales'], q=3, labels=['Low', 'Medium', 'High'])
print(df)
```

# 🏆 4. Building and Evaluating Simple Machine Learning Pipelines

Combining data preprocessing, feature engineering, and modeling into a single pipeline ensures **reproducibility** and **efficiency**.

## 🔳 4.1 Example: End-to-End Pipeline for Regression

```python
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_squared_error

# Sample housing dataset
housing = pd.DataFrame({
    'sqft': [1000, 1500, 2000, 2500, 3000],
    'bedrooms': [2, 3, 3, 4, 4],
    'price': [200000, 250000, 300000, 400000, 450000]
})

# Features and target
X = housing[['sqft', 'bedrooms']]
y = housing['price']

# Build pipeline
pipeline = Pipeline([
    ('scaler', StandardScaler()),
    ('model', RandomForestRegressor(n_estimators=100, random_state=42))
])

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Train model
pipeline.fit(X_train, y_train)

# Predictions and evaluation
y_pred = pipeline.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
print(f"Mean Squared Error: {mse:.2f}")
```

## ⚡ 4.2 Hyperparameter Tuning with GridSearchCV

Optimize model performance using **GridSearchCV**.

```python
from sklearn.model_selection import GridSearchCV

# Define parameter grid
param_grid = {
    'model__n_estimators': [50, 100, 200],
    'model__max_depth': [None, 5, 10]
}

# Apply GridSearchCV
grid_search = GridSearchCV(pipeline, param_grid, cv=3)
grid_search.fit(X_train, y_train)
```

```
# Best parameters
print("Best Parameters:", grid_search.best_params_)
```

---

# ☑ **Chapter Summary**

◈ Preprocessed and engineered features for ML models ◈ Built and evaluated pipelines using **Scikit-learn**
◈ Applied **feature engineering** with Pandas and Numpy ◈ Implemented hyperparameter tuning for model optimization

---

```
# Best parameters
print("Best Parameters:", grid_search.best_params_)
```

# 📖 Chapter 22: Best Practices and Next Steps

As you conclude your journey through data analysis with Python, it's crucial to establish best practices that ensure **clean, efficient, and scalable code**. In this chapter, we'll focus on coding standards, debugging techniques, memory optimization strategies, and how to build an impressive data analysis portfolio.

By the end of this chapter, you will be able to:

☑ Write optimized and maintainable data analysis code ☑ Debug, profile, and optimize memory usage in your data pipelines ☑ Maintain data quality and consistency across projects ☑ Build a professional data analysis portfolio to showcase your skills

## 👦💻 1. Writing Clean, Efficient, and Scalable Code

Efficient code isn't just about speed; it's about readability, maintainability, and scalability. Following clean coding principles allows you to build robust data pipelines that can evolve with project needs.

### ◍ 1.1 Code Structuring and Modularity

- **Use Functions and Classes:** Encapsulate logic in reusable functions or classes.
- **Follow PEP8 Guidelines:** Stick to Python's standard style guide for readability.
- **Leverage Type Hints:** Use type annotations for clarity.

📎 **Example: Modular Data Cleaning Function**

```python
import pandas as pd
from typing import List

def clean_dataframe(df: pd.DataFrame, columns_to_drop: List[str]) -> pd.DataFrame:
    """
    Cleans the dataframe by dropping specified columns and handling missing
values.

    Args:
        df (pd.DataFrame): Input DataFrame.
        columns_to_drop (List[str]): List of columns to drop.

    Returns:
        pd.DataFrame: Cleaned DataFrame.
    """
    df = df.drop(columns=columns_to_drop)
    df = df.dropna()
    return df

# Usage
data = {'name': ['Alice', 'Bob', None], 'age': [25, None, 30], 'city': ['NY',
'LA', 'SF']}
```

```python
df = pd.DataFrame(data)
cleaned_df = clean_dataframe(df, ['city'])
print(cleaned_df)
```

## ⚡ 1.2 Writing Vectorized and Pythonic Code

Avoid explicit Python loops in data-heavy tasks. Use **vectorized operations** for speed and simplicity.

```python
import numpy as np

# Inefficient loop-based sum
arr = np.arange(1_000_000)
total = 0
for num in arr:
    total += num

# Vectorized sum
total_vectorized = np.sum(arr)
print(f"Vectorized Sum: {total_vectorized}")
```

💡 **Tip:** Numpy and Pandas operations are optimized in C, making them significantly faster than native Python loops.

## ◍ 1.3 Scaling Code for Large Datasets

- **Use Generators:** For streaming large datasets.
- **Batch Processing:** Divide data into chunks using Dask or Vaex.
- **Out-of-Core Computation:** Use libraries like Dask for datasets that don't fit into memory.

```python
import pandas as pd

# Reading large CSV in chunks
chunk_iter = pd.read_csv('large_dataset.csv', chunksize=10000)
for chunk in chunk_iter:
    print(chunk.mean())
```

# 🐛 2. Debugging, Profiling, and Memory Optimization

Even the cleanest code can have bottlenecks. Efficient debugging and profiling can uncover hidden inefficiencies.

## ◍ 2.1 Debugging Techniques

- **pdb (Python Debugger):** Set breakpoints and step through code.
- **Logging Over Print:** Use the `logging` library for better traceability.

```python
import logging

logging.basicConfig(level=logging.INFO)
logger = logging.getLogger(__name__)

def divide(a, b):
    try:
        return a / b
    except ZeroDivisionError as e:
        logger.error("Attempted division by zero")
        return None


divide(10, 0)
```

## ⚡ 2.2 Profiling Code Performance

Use profiling tools to pinpoint slow code sections:

- `cProfile` – Built-in profiler for Python.
- `line_profiler` – Profile specific functions line-by-line.
- `memory_profiler` – Track memory usage.

```python
import cProfile

def slow_function():
    total = sum([i for i in range(1_000_000)])
    return total

cProfile.run('slow_function()')
```

## ◍ 2.3 Memory Optimization Tips

- Use **categorical data types** in Pandas to reduce memory footprint.
- Apply **vectorized operations** instead of loops.
- Process large datasets using **Dask** or **Vaex**.

```python
import pandas as pd
import numpy as np

# Optimize memory usage
df = pd.DataFrame({
    'category': ['A', 'B', 'C', 'A', 'B'],
```

```
        'value': np.random.randint(1, 100, size=5)
})

df['category'] = df['category'].astype('category')
print(df.info())
```

# 📊 3. Ensuring Data Quality and Consistency

Data quality is the backbone of any meaningful analysis. Inconsistent or inaccurate data can derail your insights.

## ◓ 3.1 Establishing Data Validation Rules

- Check for **missing values** and **outliers**.
- Validate data types and enforce schema consistency.
- Apply data integrity checks using tools like **pandas_schema** or **Great Expectations**.

```
from pandas_schema import Column, Schema
from pandas_schema.validation import CanConvertValidation, InRangeValidation

schema = Schema([
    Column('age', [CanConvertValidation(int), InRangeValidation(0, 120)]),
    Column('income', [CanConvertValidation(float)])
])

errors = schema.validate(pd.DataFrame({'age': [25, -5, 130], 'income': ['50000',
'60000', 'invalid']}))
for error in errors:
    print(error)
```

## ⚡ 3.2 Ensuring Reproducibility

- Use **random seeds** in all stochastic operations.
- Store all transformations and cleaning steps in scripts, not notebooks.
- Maintain a clear data lineage: track every change to the data.

# 📁 4. Building Your Data Analysis Portfolio

A strong portfolio showcases your expertise and practical skills. Here's how to make yours stand out.

## ◓ 4.1 What to Include in Your Portfolio?

1. **End-to-End Projects:** Show the complete data analysis pipeline from raw data to insights.
2. **Diverse Domains:** Finance, healthcare, social media—show versatility.

3. **Interactive Visualizations:** Use **Plotly** or **Dash** to build engaging dashboards.

4. **Technical Write-Ups:** Share your thought process, challenges, and insights.

---

## 📊 4.2 Portfolio Project Ideas

| Project | Skills Demonstrated |
|---|---|
| E-commerce Sales Analysis | Data cleaning, EDA, visualization |
| Stock Market Trend Prediction | Time series analysis, feature engineering |
| Social Media Sentiment Analysis | NLP, text data wrangling, visual storytelling |
| Real Estate Price Prediction | Regression modeling, feature engineering |
| COVID-19 Data Dashboard | Real-time analytics, dashboard creation |

## 🌐 4.3 Sharing Your Work

- **GitHub:** Host code repositories and notebooks.
- **Kaggle:** Participate in competitions and share notebooks.
- **Medium/Blogs:** Write about your projects and data findings.
- **LinkedIn:** Share portfolio highlights and engage with the data community.

---

# ☑️ Chapter Summary

◈ Write clean, scalable, and efficient code ◈ Debug, profile, and optimize memory usage in data pipelines ◈ Implement data validation rules for consistency and accuracy ◈ Build a diverse and engaging data analysis portfolio

---

# 📚 Appendices

The appendices provide additional resources, cheat sheets, and practical tips to solidify your understanding and improve your workflow in data analysis. Whether you're revising key concepts or looking for resources to dive deeper, this section will serve as a valuable reference.

## A A. Numpy, Pandas, and Matplotlib Cheat Sheets

### 📐 Numpy Cheat Sheet

| Operation | Example | Result |
|---|---|---|
| Create array | `np.array([1, 2, 3])` | `[1 2 3]` |
| Zeros/Ones array | `np.zeros((2, 3))` | `[[0. 0. 0.], [0. 0. 0.]]` |
| Random numbers | `np.random.rand(3, 2)` | Random 3x2 array |
| Element-wise operations | `arr * 2` | Multiplies each element by 2 |
| Matrix multiplication | `np.dot(A, B)` | Matrix product of A and B |
| Reshape array | `arr.reshape((3, 2))` | New 3x2 shaped array |
| Indexing/Slicing | `arr[1:3]` | Selects elements 1 and 2 |
| Boolean indexing | `arr[arr > 5]` | Elements greater than 5 |
| Aggregate functions | `np.mean(arr)` | Mean of array elements |

### 📊 Pandas Cheat Sheet

| Operation | Example | Result |
|---|---|---|
| Create DataFrame | `pd.DataFrame(data)` | Creates DataFrame from dict |
| Read CSV | `pd.read_csv('data.csv')` | Loads CSV into DataFrame |
| Select columns | `df['column_name']` | Series of the selected column |
| Filter rows | `df[df['Age'] > 30]` | Rows where Age > 30 |
| Group by and aggregate | `df.groupby('Dept').sum()` | Aggregated data by department |
| Apply function | `df['col'].apply(lambda x: x*2)` | Applies lambda function |
| Merge DataFrames | `pd.merge(df1, df2, on='ID')` | Joins two DataFrames on 'ID' |
| Handle missing data | `df.dropna()` / `df.fillna(0)` | Drop or fill NaN values |
| Pivot tables | `df.pivot_table(index='Dept')` | Summarizes data by department |
| Save DataFrame to CSV | `df.to_csv('output.csv')` | Writes DataFrame to CSV file |

## ☑ Matplotlib Cheat Sheet

```python
import matplotlib.pyplot as plt
import numpy as np

# Sample data
x = np.linspace(0, 10, 100)
y = np.sin(x)

# Basic Plot
plt.plot(x, y, label='Sine Wave', color='blue')
plt.title("Sine Wave Plot")
plt.xlabel("X-axis")
plt.ylabel("Y-axis")
plt.legend()
plt.show()
```

| Plot Type | Function |
|-----------|----------|
| Line Plot | `plt.plot()` |
| Bar Chart | `plt.bar()` |
| Histogram | `plt.hist()` |
| Scatter Plot | `plt.scatter()` |
| Heatmap (Seaborn) | `sns.heatmap()` |
| Subplots | `plt.subplot()` |
| Multiple Lines | `plt.plot(x, y1); plt.plot(x, y2)` |

# B. Performance Optimization Tips for Data Analysts

### ⚡ 1. Numpy and Pandas Optimization

- **Use Vectorized Operations:** Avoid Python loops where possible.
- **Minimize Memory Usage:** Use `astype('category')` for categorical data and `downcast` numeric columns.
- **Chunk Large Files:** Read large CSV files using `chunksize` in `pd.read_csv()`.

### ▦ 2. Code Profiling and Benchmarking

- **Use `cProfile` and `line_profiler`** for detailed performance insights.
- **Leverage `memory_profiler`** to track memory bottlenecks.

```python
%load_ext memory_profiler

@profile
```

```python
def memory_intensive_func():
    import numpy as np
    arr = np.random.rand(10000000)
    return np.mean(arr)
```

## 🚀 3. Parallel and Distributed Computing

- **Use Dask or Vaex** for out-of-core DataFrame manipulation.
- **Apply joblib or multiprocessing** for parallel computations.

## 🔘 C. Common Data Analysis Pitfalls and How to Avoid Them

| Pitfall | How to Avoid |
|---------|-------------|
| ⬢ Ignoring Missing Data | Use `.isnull()`, `.dropna()`, or `.fillna()` |
| ⬢ Overfitting in EDA | Stick to EDA's exploratory nature before modeling |
| ⬢ Incorrect Data Types | Check with `df.dtypes` and convert when needed |
| ⬢ Misinterpreting Correlation | Correlation ≠ Causation. Use domain knowledge. |
| ⬢ Inefficient Joins and Merges | Use indexed joins and avoid unnecessary full merges |
| ⬢ Memory Errors with Large Files | Apply chunking, Dask, or reduce DataFrame size |
| ⬢ Ignoring Time Zones in Time Series | Always localize timestamps using `.tz_localize()` |

## 🅂 D. Sample Datasets for Practice and Case Studies

| Dataset | Source | Domain |
|---------|--------|--------|
| Titanic Survival Data | Kaggle | Classification |
| NYC Taxi Trips | NYC Open Data | Time Series/Geo |
| Netflix Movies & Shows | Kaggle | EDA/Visualization |
| E-commerce Purchases | UCI Machine Learning Repository | Customer Behavior |
| Financial Market Data | Yahoo Finance API | Financial Analysis |
| COVID-19 Dataset | Johns Hopkins University | Time Series Analysis |
| Airbnb Listings | Inside Airbnb | Price Prediction |

## 📖 E. Resources for Further Learning (Courses, Books, and Tools)

### 📚 Books

- **"Python for Data Analysis"** – Wes McKinney *(The creator of Pandas)*
- **"Data Science from Scratch"** – Joel Grus

- **"Effective Pandas"** – Matt Harrison
- **"Storytelling with Data"** – Cole Nussbaumer Knaflic

---

## 🎓 Online Courses

- **Data Analysis with Python (Coursera)** – University of Michigan
- **Python for Data Science (Udemy)**
- **Practical Data Science (DataCamp)**

---

## 🛠️ Useful Tools & Libraries

- **Jupyter Notebooks** – Interactive coding and visualization
- **Seaborn** – Statistical data visualization
- **Plotly/Dash** – Interactive dashboards
- **Dask & Vaex** – Big Data handling with Pandas-like syntax
- **Great Expectations** – Data validation framework

---

# ✅ Final Thoughts

Congratulations on completing this book! 🎉 You've mastered the tools and techniques needed to perform high-level data analysis with Python.

As your next steps:

- **Explore advanced domains** like Machine Learning or Data Engineering
- **Contribute to open-source data projects**
- **Build a strong portfolio** and share your work with the data community

*"The goal is to turn data into information, and information into insight."* — Carly Fiorina

Happy Analyzing! 🚀