# PANDAS

## A Practical Guide for Data Analysts

Adriano Venancio

# About This Book

**"Pandas: A Practical Guide for Data Analysts"** by **Adriano Venancio** is designed for those who prefer to learn by doing. This book is not your typical theory-heavy resource—it's a hands-on guide created for people who want to practice coding while mastering the Pandas library.

With **12 well-structured chapters**, this ebook focuses on practical, real-world applications. Each chapter includes **clear explanations, step-by-step instructions, and code examples** that you can follow and apply directly. From the fundamentals to advanced techniques like time series analysis, performance optimization, and data visualization, this guide empowers you to solve complex data challenges with confidence.

Whether you're a beginner looking to build a strong foundation or an experienced analyst seeking advanced insights, this ebook will help you make the most of Pandas and its integration with Python libraries like Matplotlib and Seaborn.

By the end of this book, you'll have the skills to tackle real-world data problems and enhance your data analysis capabilities through practical coding experience.

# SUMMARY

# Chapter 1: Introduction to Pandas

## What is Pandas?

Pandas is an open-source Python library designed for data manipulation and analysis. Built on top of NumPy, it provides powerful, flexible, and easy-to-use data structures, such as Series and DataFrames, which allow for efficient handling of structured data. Pandas is widely used in data science, analytics, and machine learning workflows, making it a must-know tool for data professionals.

Key features of Pandas include:

- Handling of missing data.
- Alignment and reshaping of data.
- Grouping and aggregation.
- Integration with other libraries, such as Matplotlib and NumPy.

## Installation and Setup

To get started with Pandas, you need to install it. Pandas requires Python (3.7 or higher) and works well with popular distributions such as Anaconda.

### Installing Pandas with pip:

```
pip install pandas
```

### Installing Pandas with Anaconda:

If you are using the Anaconda distribution, Pandas is typically pre-installed. You can update it using:

```
conda update pandas
```

### Verifying Installation:

To check if Pandas is installed correctly, run the following in your Python environment:

```python
import pandas as pd
print(pd.__version__)
```

This will print the installed version of Pandas.

# Understanding Series and DataFrames

Pandas is built around two primary data structures:

## Series

A Series is a one-dimensional labeled array capable of holding any data type (e.g., integers, strings, floats). The labels, known as the index, allow for fast data access.

**Creating a Series:**

```python
import pandas as pd

data = [1, 2, 3, 4]
series = pd.Series(data)
print(series)
```

**Output:**

```
0    1
1    2
2    3
3    4
dtype: int64
```

The left column represents the index, and the right column represents the data.

# DataFrame

A DataFrame is a two-dimensional labeled data structure, similar to a table in a relational database or an Excel spreadsheet. It consists of rows and columns, where each column can hold different data types.

**Creating a DataFrame:**

```python
import pandas as pd

data = {
    'Name': ['Alice', 'Bob', 'Charlie'],
    'Age': [25, 30, 35],
    'City': ['New York', 'Los Angeles', 'Chicago']
}

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

**Output:**

```
      Name  Age         City
0    Alice   25     New York
1      Bob   30  Los Angeles
2  Charlie   35      Chicago
```

# Key Operations with Series and DataFrames

- **Accessing Data:**
  - Series: `series[1]` returns the value at index 1.
  - DataFrame: `df['Name']` returns the column named 'Name'.
- **Basic Descriptions:** Use `df.info()` and `df.describe()` to understand the structure and summary statistics of your data.

# Chapter 2: Pandas Basics

## Creating Series and DataFrames

Pandas provides intuitive methods to create its core data structures: Series and DataFrames.

### Creating a Series

A Series can be created from lists, dictionaries, or scalar values. Examples include:

**From a List:**

```python
import pandas as pd

numbers = [10, 20, 30, 40]
series = pd.Series(numbers)
print(series)
```

**Output:**

```
0    10
1    20
2    30
3    40
dtype: int64
```

**From a Dictionary:**

```python
data = {'a': 100, 'b': 200, 'c': 300}
series = pd.Series(data)
print(series)
```

**Output:**

```
a    100
b    200
c    300
```

```
dtype: int64
```

## Creating a DataFrame

DataFrames can be constructed from lists of dictionaries, dictionaries of lists, or external files.

**From a Dictionary of Lists:**

```
data = {
    'Name': ['Alice', 'Bob', 'Charlie'],
    'Age': [25, 30, 35],
    'City': ['New York', 'Los Angeles', 'Chicago']
}
df = pd.DataFrame(data)
print(df)
```

**Output:**

```
      Name  Age         City
0    Alice   25     New York
1      Bob   30  Los Angeles
2  Charlie   35      Chicago
```

**From a CSV File:**

```
df = pd.read_csv('data.csv')
print(df.head())
```

This loads the first few rows of the CSV file named `data.csv`.

# Indexing and Selecting Data

Pandas offers multiple ways to access and manipulate data, including label-based and position-based indexing.

## Accessing Data with `.loc[]` and `.iloc[]`

**Using `.loc[]` for Label-Based Indexing:**

```python
# Example DataFrame
data = {
    'Name': ['Alice', 'Bob', 'Charlie'],
    'Age': [25, 30, 35]
}
df = pd.DataFrame(data, index=['A', 'B', 'C'])

# Accessing rows by label
a = df.loc['A']
print(a)
```

**Output:**

```
Name     Alice
Age         25
Name: A, dtype: object
```

**Using `.iloc[]` for Position-Based Indexing:**

```python
# Accessing the first row
a = df.iloc[0]
print(a)
```

**Output:**

```
Name     Alice
Age         25
Name: A, dtype: object
```

## Accessing Columns

Columns can be accessed using the column name:

```
ages = df['Age']
print(ages)
```

**Output:**

```
A    25
B    30
C    35
Name: Age, dtype: int64
```

# Common Attributes and Methods

## Attributes:

- `df.shape`:

  ```
  import pandas as pd

  data = {
      'Name': ['Alice', 'Bob', 'Charlie'],
      'Age': [25, 30, 35],
      'City': ['New York', 'Los Angeles', 'Chicago']
  }
  df = pd.DataFrame(data)

  print(df.shape)
  ```

  **Output:**

  ```
  (3, 3)
  ```

  This indicates 3 rows and 3 columns.

- `df.columns`:

```
print(df.columns)
```

**Output:**

```
Index(['Name', 'Age', 'City'], dtype='object')
```

- `df.index`:

```
print(df.index)
```

**Output:**

```
RangeIndex(start=0, stop=3, step=1)
```

## Methods:

- `df.head(n)`:

```
print(df.head(2))
```

**Output:**

```
    Name  Age         City
0  Alice   25     New York
1    Bob   30  Los Angeles
```

- `df.tail(n)`:

```
print(df.tail(1))
```

**Output:**

```
      Name  Age     City
2  Charlie   35  Chicago
```

- `df.info()`:

```
print(df.info())
```

**Output:**

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3 entries, 0 to 2
Data columns (total 3 columns):
 #   Column  Non-Null Count  Dtype
---  ------  --------------  -----
 0   Name    3 non-null      object
 1   Age     3 non-null      int64
 2   City    3 non-null      object
dtypes: int64(1), object(2)
memory usage: 200.0+ bytes
```

- `df.describe()`:

```
print(df.describe())
```

**Output:**

```
        Age
count   3.0
mean   30.0
std     5.0
min    25.0
25%    27.5
50%    30.0
75%    32.5
```

```
max      35.0
```

# Chapter 3: Data Manipulation with Pandas

Data manipulation is one of the most powerful features of Pandas. This chapter explores how to transform, filter, and organize your data effectively.

## Adding and Dropping Columns/Rows

### Adding a New Column:

New columns can be added by assigning values to a column name.

```python
import pandas as pd

data = {
    'Name': ['Alice', 'Bob', 'Charlie'],
    'Age': [25, 30, 35]
}
df = pd.DataFrame(data)

# Add a new column

df['Salary'] = [50000, 60000, 70000]
print(df)
```

**Output:**

```
      Name  Age  Salary
0    Alice   25   50000
1      Bob   30   60000
2  Charlie   35   70000
```

## Dropping Columns or Rows:

To drop columns or rows, use the `drop()` method.

**Dropping a Column:**

```python
# Drop the 'Salary' column
df = df.drop('Salary', axis=1)
print(df)
```

**Output:**

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

**Dropping a Row:**

```python
# Drop the row at index 1
df = df.drop(1, axis=0)
print(df)
```

**Output:**

```
      Name  Age
0    Alice   25
2  Charlie   35
```

# Renaming Columns and Indexes

To rename columns or indexes, use the `rename()` method.

## Renaming Columns:

```python
df = df.rename(columns={'Name': 'Full Name'})
print(df)
```

**Output:**

```
     Full Name  Age
0     Alice      25
1       Bob      30
2   Charlie      35
```

## Renaming Indexes:

```python
df = df.rename(index={0: 'Row1', 1: 'Row2'})
print(df)
```

**Output:**

```
         Full Name  Age
Row1       Alice     25
Row2         Bob     30
2        Charlie     35
```

# Filtering Data

Filtering helps to select subsets of data based on conditions.

## Filtering Rows Based on a Condition:

```python
# Filter rows where Age > 25
filtered_df = df[df['Age'] > 25]
print(filtered_df)
```

**Output:**

```
      Name  Age
1      Bob   30
2  Charlie   35
```

## Filtering with Multiple Conditions:

```python
# Filter rows where Age > 25 and Name starts with 'C'
filtered_df = df[(df['Age'] > 25) & (df['Name'].str.startswith('C'))]
print(filtered_df)
```

**Output:**

```
      Name  Age
2  Charlie   35
```

# Sorting Data

Sorting is essential for organizing data for better readability or further analysis.

## Sorting by Column Values:

```python
# Sort by Age in ascending order
df = df.sort_values(by='Age')
print(df)
```

**Output:**

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

## Sorting in Descending Order:

```python
# Sort by Age in descending order
df = df.sort_values(by='Age', ascending=False)
print(df)
```

**Output:**

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

# Chapter 4: Grouping and Aggregation in Pandas

Grouping and aggregation allow you to analyze and summarize data effectively. In this chapter, we will explore how to use `groupby()` and other aggregation functions to gain insights from your data.

## Grouping Data with `groupby()`

The `groupby()` method is used to split data into groups based on some criteria. Once grouped, you can apply aggregation functions to summarize the data.

### Basic Grouping

```python
import pandas as pd

data = {
    'Department': ['HR', 'HR', 'IT', 'IT', 'Finance'],
    'Employee': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
    'Salary': [50000, 60000, 75000, 80000, 70000]
}
df = pd.DataFrame(data)

# Group by Department
grouped = df.groupby('Department')

# Calculate the mean salary by department
mean_salary = grouped['Salary'].mean()
print(mean_salary)
```

**Output:**

```
Department
Finance    70000.0
HR         55000.0
IT         77500.0
Name: Salary, dtype: float64
```

# Applying Multiple Aggregations

You can use multiple aggregation functions simultaneously using the `agg()` method.

## Example:

```python
# Aggregate Salary with multiple functions
grouped_agg = grouped['Salary'].agg(['mean', 'max', 'min'])
print(grouped_agg)
```

**Output:**

```
              mean      max      min
Department
Finance     70000.0   70000   70000
HR          55000.0   60000   50000
IT          77500.0   80000   75000
```

# Filtering Groups

Sometimes, you may want to filter groups based on specific conditions.

## Example:

```python
# Filter groups where the mean salary is greater than 60000
filtered = grouped.filter(lambda x: x['Salary'].mean() > 60000)
print(filtered)
```

**Output:**

```
  Department Employee  Salary
2         IT  Charlie   75000
3         IT    David   80000
4    Finance      Eve   70000
```

# Pivot Tables

Pivot tables are a powerful way to summarize data, similar to Excel pivot tables.

## Example:

```python
# Create a pivot table to summarize salaries by Department
pivot = df.pivot_table(values='Salary', index='Department', aggfunc='mean')
print(pivot)
```

**Output:**

```
             Salary
Department
Finance     70000.0
HR          55000.0
IT          77500.0
```

# Custom Aggregation Functions

You can create custom functions to perform specific aggregations.

## Example:

```python
# Define a custom function to calculate the salary range
def salary_range(series):
    return series.max() - series.min()

# Apply the custom function
grouped_custom = grouped['Salary'].agg(salary_range)
print(grouped_custom)
```

**Output:**

```
Department
Finance         0
HR          10000
IT           5000
Name: Salary, dtype: int64
```

# Chapter 5: Cleaning Data with Pandas

Data cleaning is a crucial step in any data analysis or machine learning project. Pandas, a powerful Python library for data manipulation and analysis, provides a plethora of tools to clean and preprocess your data efficiently. In this chapter, we will explore the most effective techniques and best practices for cleaning data using Pandas.

## Why is Data Cleaning Important?

Data cleaning is essential to ensure the accuracy and reliability of your analysis. Raw data often contains:

- Missing values
- Duplicates
- Outliers
- Incorrect or inconsistent data types
- Irrelevant information

Cleaning data helps you uncover valuable insights and build robust models.

Load your dataset into a DataFrame:

```python
import pandas as pd

data = pd.DataFrame({
    'Name': ['Alice', 'Bob', 'Charlie', None],
    'Age': [25, None, 30, 22],
    'Salary': [50000, 60000, None, 40000]
})
print(data)
```

**Output:**

```
      Name   Age   Salary
0    Alice  25.0  50000.0
1      Bob   NaN  60000.0
```

```
2  Charlie  30.0      NaN
3     None  22.0  40000.0
```

# 5.1. Handling Missing Data

## Identifying Missing Values

Use the `isnull()` method to identify missing values:

```
missing_data = data.isnull()
print(missing_data)
print(missing_data.sum())
```

**Output:**

```
    Name    Age  Salary
0  False  False   False
1  False   True   False
2  False  False    True
3   True  False   False
Name      1
Age       1
Salary    1
dtype: int64
```

## Filling Missing Values

- **Fill with a specific value:**

```
data['Age'].fillna(0, inplace=True)
print(data)
```

**Output:**

```
        Name   Age    Salary
0      Alice  25.0  50000.0
1        Bob   0.0  60000.0
2    Charlie  30.0      NaN
3       None  22.0  40000.0
```

- **Fill with the mean/median/mode:**

```
mean_salary = data['Salary'].mean()
data['Salary'].fillna(mean_salary, inplace=True)
print(data)
```

**Output:**

```
        Name   Age    Salary
0      Alice  25.0  50000.0
1        Bob   0.0  60000.0
2    Charlie  30.0  50000.0
3       None  22.0  40000.0
```

## Dropping Missing Values

- **Drop rows with missing values:**

```
data = pd.DataFrame({
    'Name': ['Alice', 'Bob', 'Charlie', None],
    'Age': [25, None, 30, 22],
    'Salary': [50000, 60000, None, 40000]
})
data.dropna(inplace=True)
print(data)
```

**Output:**

```
      Name   Age    Salary
0    Alice  25.0   50000.0
```

- **Drop columns with missing values:**

```python
data = pd.DataFrame({
    'Name': ['Alice', 'Bob', 'Charlie', None],
    'Age': [25, None, 30, 22],
    'Salary': [50000, 60000, None, 40000]
})
data.dropna(axis=1, inplace=True)
print(data)
```

**Output:**

```
      Name
0    Alice
1      Bob
2  Charlie
3     None
```

# 5.2. Removing Duplicates

## Identifying Duplicates

```python
data = pd.DataFrame({
    'Name': ['Alice', 'Alice', 'Bob', 'Charlie'],
    'Age': [25, 25, 22, 30]
})
duplicates = data.duplicated()
print(duplicates)
```

**Output:**

```
0    False
1     True
2    False
3    False
dtype: bool
```

## Dropping Duplicates

```python
data.drop_duplicates(inplace=True)
print(data)
```

**Output:**

```
      Name  Age
0     Alice   25
2       Bob   22
3   Charlie   30
```

# 5.3. Handling Outliers

## Detecting Outliers

Use statistical methods like the Interquartile Range (IQR):

```python
data = pd.DataFrame({'Age': [22, 25, 30, 120]})
Q1 = data['Age'].quantile(0.25)
Q3 = data['Age'].quantile(0.75)
IQR = Q3 - Q1

outliers = data[(data['Age'] < Q1 - 1.5 * IQR) | (data['Age'] > Q3 + 1.5 * IQR)]
print(outliers)
```

**Output:**

```
     Age
3   120
```

## Handling Outliers

- **Remove outliers:**

```python
data = data[(data['Age'] >= Q1 - 1.5 * IQR) & (data['Age'] <= Q3 + 1.5 * IQR)]
print(data)
```

**Output:**

```
    Age
0    22
1    25
2    30
```

- **Cap or floor extreme values:**

```python
data = pd.DataFrame({'Age': [22, 25, 30, 120]})
data['Age'] = data['Age'].clip(lower=Q1 - 1.5 * IQR, upper=Q3 + 1.5 * IQR)
print(data)
```

**Output:**

```
    Age
0    22
1    25
2    30
3    30
```

# 5.4. Transforming Data

## Changing Data Types

```python
data = pd.DataFrame({'Age': [25.0, 30.0, 22.0]})
data['Age'] = data['Age'].astype('int')
print(data)
```

**Output:**

```
    Age
0   25
1   30
2   22
```

## Renaming Columns

```python
data = pd.DataFrame({'old_name': [1, 2, 3]})
data.rename(columns={'old_name': 'new_name'}, inplace=True)
print(data)
```

**Output:**

```
    new_name
0          1
1          2
2          3
```

## Normalizing or Scaling Data

```python
data = pd.DataFrame({'Age': [22, 25, 30]})
data['Age'] = (data['Age'] - data['Age'].mean()) / data['Age'].std()
print(data)
```

**Output:**

```
        Age
0 -1.135550
1 -0.162221
2  1.297771
```

## Encoding Categorical Data

- **One-hot encoding:**

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

**Output:**

```
   Category_A  Category_B
0           1           0
1           0           1
2           1           0
```

- **Label encoding:**

```
data = pd.DataFrame({'Category': ['A', 'B', 'A']})
data['Category'] = data['Category'].astype('category').cat.codes
print(data)
```

**Output:**

```
   Category
0         0
```

```
1        1
2        0
```

## 5.5. Working with Strings

### Removing Unnecessary Characters

```python
data = pd.DataFrame({'Text': ['Hello!', 'World@', '#Python']})
data['Text'] = data['Text'].str.replace(r'[^a-zA-Z0-9]', '', regex=True)
print(data)
```

**Output:**

```
     Text
0   Hello
1   World
2  Python
```

### Converting Case

- **To lowercase:**

```python
data['Text'] = data['Text'].str.lower()
print(data)
```

**Output:**

```
     Text
0   hello
1   world
2  python
```

- **To uppercase:**

```python
data['Text'] = data['Text'].str.upper()
print(data)
```

**Output:**

```
     Text
0   HELLO
1   WORLD
2  PYTHON
```

## Splitting and Combining Columns

- **Splitting:**

```python
data = pd.DataFrame({'full_name': ['Alice Smith', 'Bob Brown']})
data[['first', 'last']] = data['full_name'].str.split(' ', expand=True)
print(data)
```

**Output:**

```
     full_name  first    last
0  Alice Smith  Alice   Smith
1    Bob Brown    Bob   Brown
```

- **Combining:**

```python
data['full_name'] = data['first'] + ' ' + data['last']
print(data)
```

**Output:**

```
     full_name
0  Alice Smith
1   Bob Brown
```

# 5.6. Filtering and Selecting Data

## Filtering Rows

```python
data = pd.DataFrame({'Age': [22, 25, 30]})
filtered_data = data[data['Age'] > 25]
print(filtered_data)
```

**Output:**

```
   Age
2   30
```

## Selecting Specific Columns

```python
data = pd.DataFrame({'col1': [1, 2], 'col2': [3, 4]})
selected_data = data[['col1']]
print(selected_data)
```

**Output:**

```
   col1
0     1
1     2
```

## 5.7. Saving the Cleaned Data

After cleaning your data, save it for further analysis:

```python
data.to_csv('cleaned_data.csv', index=False)
print('Data saved to cleaned_data.csv')
```

**Output:**

```
Data saved to cleaned_data.csv
```

# Chapter 6: Merging, Joining, and Concatenating DataFrames

Combining datasets is an essential part of data analysis. Pandas provides various methods for merging, joining, and concatenating data.

## Merging DataFrames

The `merge()` function combines two DataFrames based on a common column or index. It is similar to SQL joins.

### Example:

```python
import pandas as pd

# DataFrame 1
data1 = {
    'ID': [1, 2, 3],
    'Name': ['Alice', 'Bob', 'Charlie']
}
df1 = pd.DataFrame(data1)

# DataFrame 2
data2 = {
    'ID': [2, 3, 4],
    'Salary': [60000, 70000, 80000]
}
df2 = pd.DataFrame(data2)

# Merge DataFrames on ID
merged_df = pd.merge(df1, df2, on='ID', how='inner')
print(merged_df)
```

**Output:**

```
   ID     Name  Salary
0   2      Bob   60000
1   3  Charlie   70000
```

## Merge Types:

- **Inner Join:** Includes only matching rows.
- **Outer Join:** Includes all rows from both DataFrames, filling with `NaN` where there is no match.
- **Left Join:** Includes all rows from the left DataFrame and matches from the right.
- **Right Join:** Includes all rows from the right DataFrame and matches from the left.

# Joining DataFrames

The `join()` method is used to combine DataFrames based on their index.

## Example:

```python
# Set index for both DataFrames
df1.set_index('ID', inplace=True)
df2.set_index('ID', inplace=True)

# Join DataFrames
joined_df = df1.join(df2, how='inner')
print(joined_df)
```

**Output:**

```
        Name  Salary
ID
2        Bob   60000
3    Charlie   70000
```

# Concatenating DataFrames

The `concat()` function stacks DataFrames either vertically (default) or horizontally.

## Vertical Concatenation:

```python
# Create two DataFrames
data1 = {
    'Name': ['Alice', 'Bob'],
    'Salary': [50000, 60000]
}
data2 = {
    'Name': ['Charlie', 'David'],
    'Salary': [70000, 80000]
}
df1 = pd.DataFrame(data1)
df2 = pd.DataFrame(data2)

# Concatenate DataFrames
concat_df = pd.concat([df1, df2])
print(concat_df)
```

**Output:**

```
      Name  Salary
0    Alice   50000
1      Bob   60000
0  Charlie   70000
1    David   80000
```

## Horizontal Concatenation:

```python
# Concatenate DataFrames along columns
concat_df_horizontal = pd.concat([df1, df2], axis=1)
print(concat_df_horizontal)
```

**Output:**

```
    Name  Salary     Name  Salary
0  Alice   50000  Charlie   70000
```

```
1       Bob    60000    David    80000
```

# Handling Duplicate Indices

When concatenating, duplicate indices can occur. Use the `ignore_index` parameter to reindex the result.

## Example:

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

**Output:**

```
      Name  Salary
0    Alice   50000
1      Bob   60000
2  Charlie   70000
3    David   80000
```

# Chapter 7: Advanced Indexing and Reshaping Data

Indexing and reshaping data are crucial for reorganizing datasets to suit specific analysis needs. In this chapter, we explore hierarchical indexing, reshaping methods, and pivoting data.

## Hierarchical Indexing

Hierarchical indexing (also known as multi-level indexing) allows you to work with data stored in a multi-dimensional manner within a DataFrame.

### Creating a MultiIndex DataFrame:

```python
import pandas as pd

# Define data
data = {
    'City': ['New York', 'Los Angeles', 'Chicago', 'Houston'],
    'Year': [2020, 2020, 2021, 2021],
    'Population': [8.3, 4.0, 2.7, 2.3]
}
df = pd.DataFrame(data)

# Set a MultiIndex
df = df.set_index(['City', 'Year'])
print(df)
```

**Output:**

```
                    Population
City        Year
New York    2020          8.3
Los Angeles 2020          4.0
Chicago     2021          2.7
Houston     2021          2.3
```

## Accessing Data with MultiIndex:

```python
# Access data for New York in 2020
ny_population = df.loc[('New York', 2020)]
print(ny_population)
```

**Output:**

```
Population    8.3
Name: (New York, 2020), dtype: float64
```

# Reshaping DataFrames

Pandas provides powerful methods to reshape data using `stack()`, `unstack()`, and `melt()`.

## Stacking and Unstacking:

### Stacking:

Converts columns into rows.

```
stacked = df.stack()
print(stacked)
```

**Output:**

```
City         Year
New York     2020  Population    8.3
Los Angeles  2020  Population    4.0
Chicago      2021  Population    2.7
Houston      2021  Population    2.3
```

### Unstacking:

Converts rows into columns.

```
unstacked = df.unstack()
print(unstacked)
```

**Output:**

```
             Population
Year           2020  2021
City
New York        8.3   NaN
Los Angeles     4.0   NaN
Chicago         NaN   2.7
Houston         NaN   2.3
```

## Melting:

Transforms a DataFrame into a long format.

```python
# Example DataFrame
data = {
    'Name': ['Alice', 'Bob'],
    'Math': [90, 80],
    'Science': [85, 95]
}
df = pd.DataFrame(data)

# Melt the DataFrame
melted = pd.melt(df, id_vars=['Name'], var_name='Subject', value_name='Score')
print(melted)
```

**Output:**

```
    Name  Subject  Score
0  Alice     Math     90
1    Bob     Math     80
2  Alice  Science     85
3    Bob  Science     95
```

# Pivoting DataFrames

The `pivot()` and `pivot_table()` methods allow you to reshape data into a wide format.

## Using `pivot()`:

```python
# Example DataFrame
data = {
    'Date': ['2023-01-01', '2023-01-02', '2023-01-01'],
    'City': ['New York', 'New York', 'Los Angeles'],
    'Temperature': [30, 32, 75]
}
df = pd.DataFrame(data)

# Pivot the DataFrame
pivoted = df.pivot(index='Date', columns='City', values='Temperature')
print(pivoted)
```

**Output:**

```
City         Los Angeles  New York
Date
2023-01-01          75.0      30.0
2023-01-02           NaN      32.0
```

## Using `pivot_table()` for Aggregation:

```python
# Pivot table with aggregation
pivot_table = df.pivot_table(values='Temperature', index='City',
aggfunc='mean')
print(pivot_table)
```

**Output:**

```
City
Los Angeles  75.0
New York     31.0
```

# Chapter 8: Working with Time Series Data in Pandas

Time series data is critical for analyzing trends and patterns over time. Pandas provides robust tools to handle, manipulate, and analyze time series data.

## Generating and Parsing Dates

Pandas allows you to work seamlessly with date and time data through the `datetime` module and its own date functions.

### Creating a Date Range:

```python
import pandas as pd

# Create a date range
date_range = pd.date_range(start='2023-01-01', end='2023-01-07')
print(date_range)
```

**Output:**

```
DatetimeIndex(['2023-01-01', '2023-01-02', '2023-01-03', '2023-01-04',
               '2023-01-05', '2023-01-06', '2023-01-07'],
              dtype='datetime64[ns]', freq='D')
```

### Converting Strings to Datetime:

```python
# Convert strings to datetime
strings = ['2023-01-01', '2023-01-02', '2023-01-03']
dates = pd.to_datetime(strings)
print(dates)
```

**Output:**

```
DatetimeIndex(['2023-01-01', '2023-01-02', '2023-01-03'],
dtype='datetime64[ns]', freq=None)
```

# Setting a Datetime Index

Datetime indices are essential for time series analysis as they allow time-based slicing and indexing.

## Example:

```python
# Example data
data = {
    'Date': ['2023-01-01', '2023-01-02', '2023-01-03'],
    'Sales': [200, 150, 300]
}
df = pd.DataFrame(data)

# Convert 'Date' column to datetime and set as index
df['Date'] = pd.to_datetime(df['Date'])
df.set_index('Date', inplace=True)
print(df)
```

**Output:**

```
            Sales
Date
2023-01-01    200
2023-01-02    150
2023-01-03    300
```

# Resampling Time Series Data

Resampling involves changing the frequency of your time series data, such as converting daily data to monthly data.

## Example:

```python
# Resample to monthly frequency
resampled = df.resample('M').sum()
print(resampled)
```

**Output:**

```
            Sales
Date
2023-01-31    650
```

# Rolling and Expanding Calculations

Rolling and expanding methods provide ways to compute statistics over a moving window or cumulative view.

## Rolling Window:

```python
# Compute rolling mean with a window of 2
df['Rolling_Mean'] = df['Sales'].rolling(window=2).mean()
print(df)
```

**Output:**

```
            Sales  Rolling_Mean
Date
2023-01-01    200           NaN
2023-01-02    150         175.0
2023-01-03    300         225.0
```

## Expanding Window:

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

**Output:**

```
            Sales  Rolling_Mean  Expanding_Mean
Date
2023-01-01    200           NaN           200.0
2023-01-02    150         175.0           175.0
2023-01-03    300         225.0           216.7
```

# Shifting and Lagging Data

Shifting data is useful for calculating differences or creating lagged features.

## Example:

```python
# Shift sales data by 1 day
df['Lagged_Sales'] = df['Sales'].shift(1)
print(df)
```

**Output:**

```
            Sales   Rolling_Mean   Expanding_Mean   Lagged_Sales
Date
2023-01-01   200            NaN            200.0            NaN
2023-01-02   150          175.0            175.0          200.0
2023-01-03   300          225.0            216.7          150.0
```

# Time Series Visualization

Visualizing time series data helps uncover trends and seasonal patterns.

## Example:

```python
import matplotlib.pyplot as plt

# Plot sales data
df['Sales'].plot(title='Sales Over Time')
plt.xlabel('Date')
plt.ylabel('Sales')
plt.show()
```

# Chapter 9: Performance Optimization in Pandas

Working with large datasets in Pandas can be challenging in terms of performance. This chapter explores techniques to optimize memory usage and improve computation speed.

## Optimizing Memory Usage

Large datasets can consume significant memory. Pandas provides options to reduce memory consumption by optimizing data types.

### Downcasting Numeric Types

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

# Example data
data = {
    'Integers': [1, 2, 3, 4, 5],
    'Floats': [1.0, 2.1, 3.2, 4.3, 5.4]
}
df = pd.DataFrame(data)

# Downcast integers and floats
df['Integers'] = pd.to_numeric(df['Integers'], downcast='integer')
df['Floats'] = pd.to_numeric(df['Floats'], downcast='float')

print(df.dtypes)
```

**Output:**

```
Integers      int8
Floats       float32
dtype: object
```

## Converting Object Columns to Categorical

```python
# Example data
data = {
    'Category': ['A', 'B', 'A', 'C', 'B']
}
df = pd.DataFrame(data)

# Convert to category
df['Category'] = df['Category'].astype('category')
print(df.dtypes)
```

**Output:**

```
Category    category
dtype: object
```

# Efficient Iterations

Avoid row-wise iterations (`apply()` or `iterrows()`) when possible. Use vectorized operations instead.

## Example of Vectorized Operation:

```python
# Example data
data = {
    'A': [1, 2, 3, 4],
    'B': [5, 6, 7, 8]
}
df = pd.DataFrame(data)

# Vectorized operation
df['Sum'] = df['A'] + df['B']
print(df)
```

**Output:**

```
   A  B  Sum
0  1  5    6
1  2  6    8
2  3  7   10
3  4  8   12
```

# Chunking Large DataFrames

When working with large datasets, reading and processing the data in chunks can prevent memory issues.

## Example:

```python
# Reading a CSV file in chunks
chunk_iter = pd.read_csv('large_file.csv', chunksize=1000)

for chunk in chunk_iter:
    # Process each chunk
    print(chunk.head())
```

# Parallel Processing with Dask

Dask extends Pandas for larger-than-memory computations by breaking data into smaller chunks and processing them in parallel.

## Example:

```python
import dask.dataframe as dd

# Create a Dask DataFrame
df = dd.read_csv('large_file.csv')

# Perform operations
df['new_column'] = df['A'] + df['B']
result = df.compute()
print(result)
```

# Avoiding Copy-on-Write Pitfalls

Minimize unnecessary copying of DataFrames to improve performance.

**Example:**

```python
# Inefficient
new_df = df.copy()
new_df['New_Column'] = new_df['A'] * 2

# Efficient
df['New_Column'] = df['A'] * 2
```

# Chapter 10: Real-World Applications of Pandas

Pandas is widely used in various real-world applications across industries. This chapter demonstrates practical use cases of Pandas for data analysis and preprocessing tasks.

## Case Study 1: Cleaning and Preprocessing Data

### Problem:

A company has sales data with missing values and inconsistent formatting. The goal is to clean the data for analysis.

### Solution:

```python
import pandas as pd

# Sample data with missing values
data = {
    'Product': ['Laptop', 'Tablet', None, 'Smartphone'],
    'Price': [1000, 500, None, 800],
    'Quantity': [10, None, 15, 20]
}
df = pd.DataFrame(data)

# Fill missing values
df['Price'] = df['Price'].fillna(df['Price'].mean())
df['Quantity'] = df['Quantity'].fillna(0)
df['Product'] = df['Product'].fillna('Unknown')

# Standardize column names
df.columns = [col.lower() for col in df.columns]

print(df)
```

**Output:**

```
      product   price  quantity
0      Laptop  1000.0      10.0
1      Tablet   500.0       0.0
2     Unknown   766.7      15.0
```

```
3   Smartphone   800.0      20.0
```

# Case Study 2: Exploratory Data Analysis (EDA)

## Problem:

Analyze a dataset to uncover insights about customer behavior.

## Solution:

```python
# Example data
data = {
    'Customer': ['Alice', 'Bob', 'Charlie', 'David'],
    'Spending': [200, 150, 300, 400],
    'Visits': [5, 3, 8, 2]
}
df = pd.DataFrame(data)

# Add derived metrics
df['Avg_Spend_per_Visit'] = df['Spending'] / df['Visits']

# Summary statistics
print(df.describe())

# Filter high-value customers
high_value_customers = df[df['Spending'] > 250]
print(high_value_customers)
```

**Output:**

```
       Spending  Visits  Avg_Spend_per_Visit
count    4.0000  4.0000               4.0000
mean   262.5000  4.5000              58.7500
std    111.8034  2.5000              27.9509
min    150.0000  2.0000              37.5000
25%    187.5000  3.2500              46.8750
50%    250.0000  4.0000              50.0000
75%    325.0000  5.2500              61.8750
```

```
max    400.0000  8.0000              100.0000


   Customer  Spending  Visits  Avg_Spend_per_Visit
2  Charlie        300       8                37.50
3    David        400       2               100.00
```

# Case Study 3: Building a Simple ETL Pipeline

## Problem:

Automate the extraction, transformation, and loading of data for analysis.

## Solution:

```python
# Extract
raw_data = {
    'Name': ['Alice', 'Bob'],
    'Age': ['25', '30'],
    'Salary': ['$5000', '$6000']
}
df = pd.DataFrame(raw_data)

# Transform
df['Age'] = pd.to_numeric(df['Age'])
df['Salary'] = df['Salary'].replace({'\$': ''}, regex=True).astype(int)

# Load (example: save to CSV)
df.to_csv('cleaned_data.csv', index=False)

print(df)
```

**Output:**

```
    Name  Age  Salary
0  Alice   25    5000
1    Bob   30    6000
```

# Chapter 11: Data Visualization with Matplotlib and Seaborn

Visualization is a key part of data analysis, and Pandas integrates seamlessly with Matplotlib and Seaborn to produce compelling visual insights. This chapter demonstrates how to leverage these libraries for effective data visualization.

## Section 1: Visualization with Matplotlib

Matplotlib is a foundational Python library for creating static, interactive, and animated visualizations.

### Example 1: Line Plot

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

# Example data
data = {
    'Month': ['Jan', 'Feb', 'Mar', 'Apr'],
    'Sales': [200, 250, 300, 400]
}
df = pd.DataFrame(data)

# Line plot
plt.plot(df['Month'], df['Sales'], marker='o', linestyle='-', color='b')
plt.title('Monthly Sales')
plt.xlabel('Month')
plt.ylabel('Sales')
plt.grid(True)
plt.show()
```

# Monthly Sales

## Example 2: Bar Plot

```python
# Bar plot
plt.bar(df['Month'], df['Sales'], color='skyblue')
plt.title('Monthly Sales (Bar Plot)')
plt.xlabel('Month')
plt.ylabel('Sales')
plt.show()
```

## Example 3: Histogram

```python
# Histogram
plt.hist(df['Sales'], bins=5, color='orange', edgecolor='black')
plt.title('Sales Distribution')
plt.xlabel('Sales')
plt.ylabel('Frequency')
plt.show()
```



Sales Distribution

## Example 4: Scatter Plot

```python
# Scatter plot
plt.scatter(df['Month'], df['Sales'], color='green')
plt.title('Monthly Sales (Scatter Plot)')
plt.xlabel('Month')
plt.ylabel('Sales')
plt.show()
```

## Example 5: Pie Chart

```python
# Pie chart
plt.pie(df['Sales'], labels=df['Month'], autopct='%1.1f%%', startangle=90)
plt.title('Sales Distribution by Month')
plt.show()
```

Sales Distribution by Month

## Example 6: Stacked Bar Plot

```python
# Stacked bar plot
data = {
    'Month': ['Jan', 'Feb', 'Mar', 'Apr'],
    'Product_A': [50, 60, 70, 80],
    'Product_B': [150, 190, 230, 320]
}
df = pd.DataFrame(data)

plt.bar(df['Month'], df['Product_A'], label='Product A', color='blue')
plt.bar(df['Month'], df['Product_B'], bottom=df['Product_A'], label='Product B', color='orange')
plt.title('Sales Distribution by Product')
plt.xlabel('Month')
plt.ylabel('Sales')
plt.legend()
plt.show()
```

## Example 7: Step Plot

```python
# Step plot
plt.step(df['Month'], df['Sales'], where='mid', color='purple')
plt.title('Monthly Sales (Step Plot)')
plt.xlabel('Month')
plt.ylabel('Sales')
plt.show()
```


Monthly Sales (Step Plot)

## Example 8: Horizontal Bar Plot

```python
# Horizontal bar plot
plt.barh(df['Month'], df['Sales'], color='cyan')
plt.title('Monthly Sales (Horizontal Bar Plot)')
plt.xlabel('Sales')
plt.ylabel('Month')
plt.show()
```



Monthly Sales (Horizontal Bar Plot)

# Section 2: Advanced Visualization with Seaborn

Seaborn builds on Matplotlib to create visually appealing and informative statistical graphics.

## Example 1: Heatmap

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

# Example data
data = np.random.rand(5, 5)
columns = ['A', 'B', 'C', 'D', 'E']
df = pd.DataFrame(data, columns=columns)

# Heatmap
sns.heatmap(df, annot=True, cmap='coolwarm', linewidths=0.5)
plt.title('Heatmap Example')
plt.show()
```
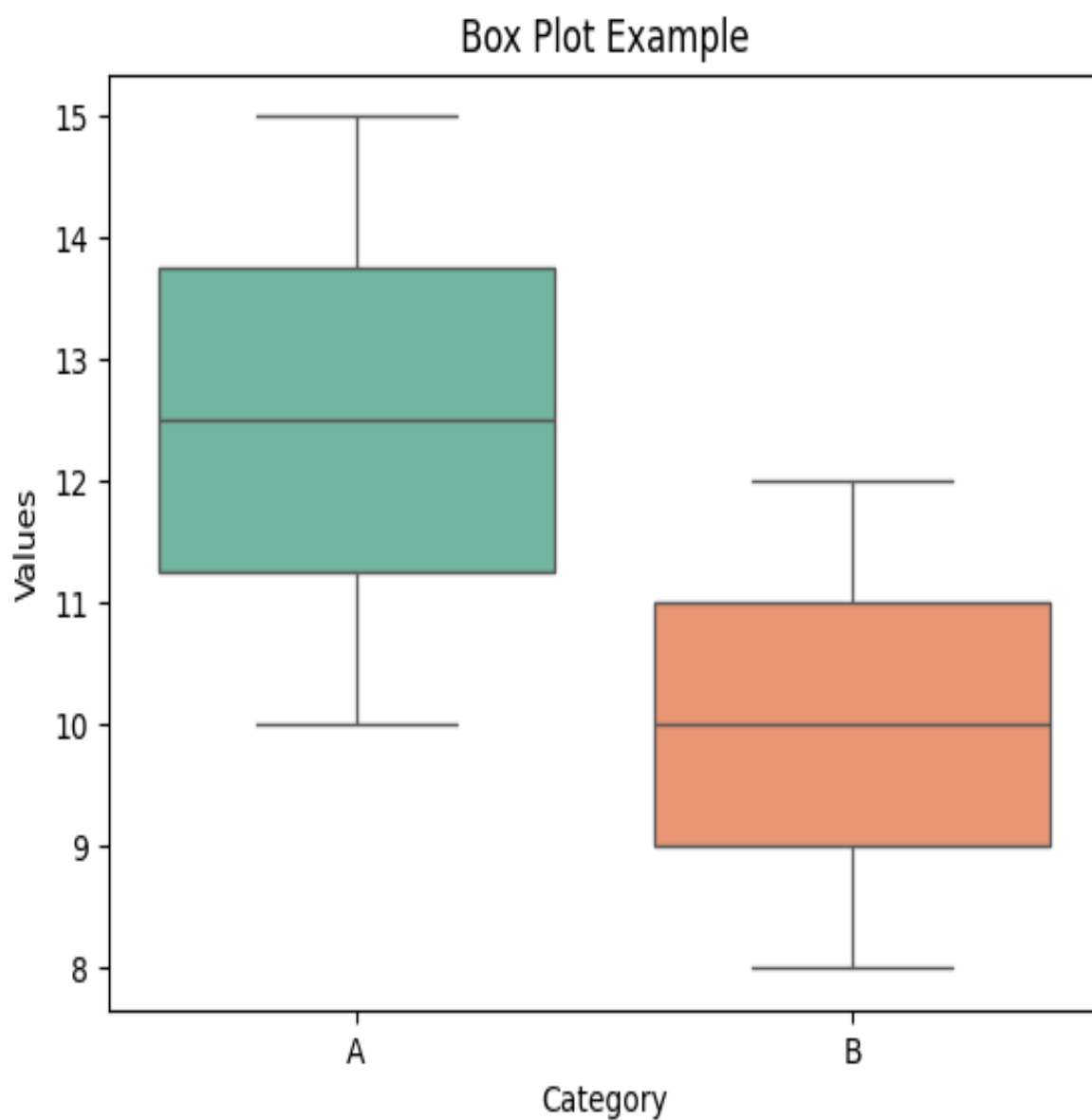
Heatmap Example

## Example 2: Box Plot

```python
# Example data
data = {
    'Category': ['A', 'A', 'B', 'B'],
    'Values': [10, 15, 8, 12]
}
df = pd.DataFrame(data)

# Box plot
sns.boxplot(x='Category', y='Values', data=df, palette='Set2')
plt.title('Box Plot Example')
plt.show()
```
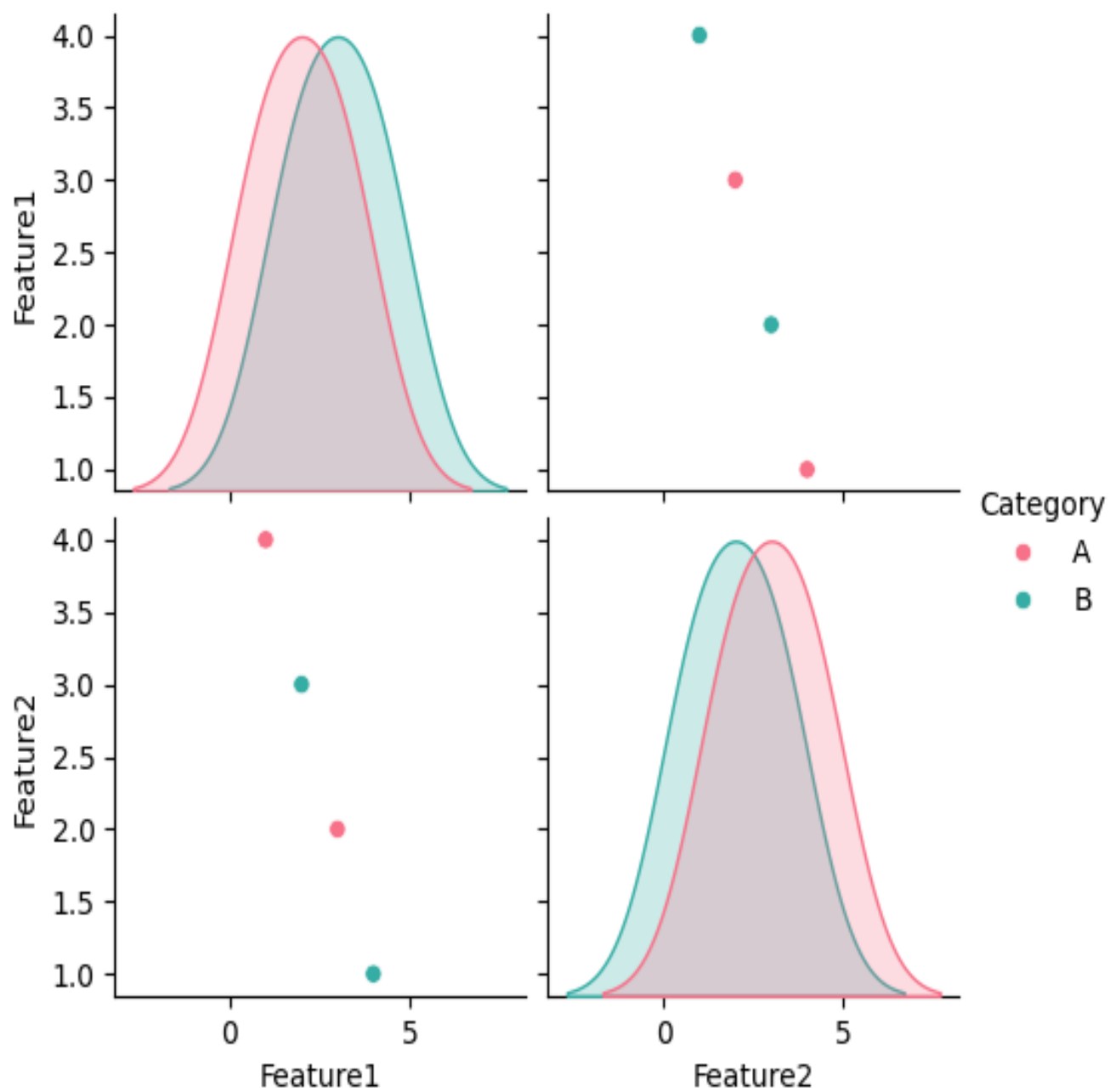


Box Plot Example

## Example 3: Pairplot

```python
# Example data
data = {
    'Feature1': [1, 2, 3, 4],
    'Feature2': [4, 3, 2, 1],
    'Category': ['A', 'B', 'A', 'B']
}
df = pd.DataFrame(data)

# Pairplot
sns.pairplot(df, hue='Category', palette='husl')
plt.show()
```
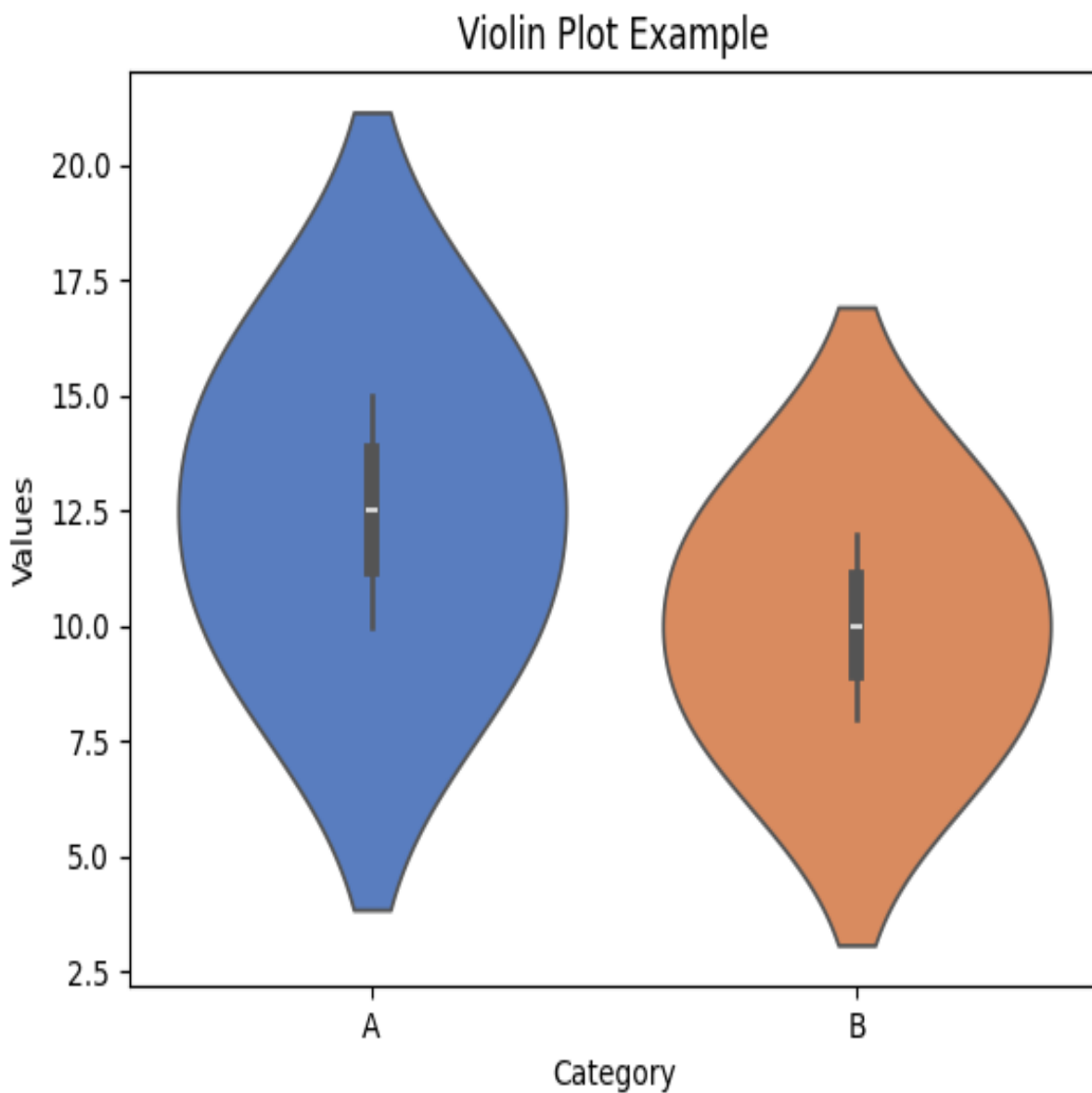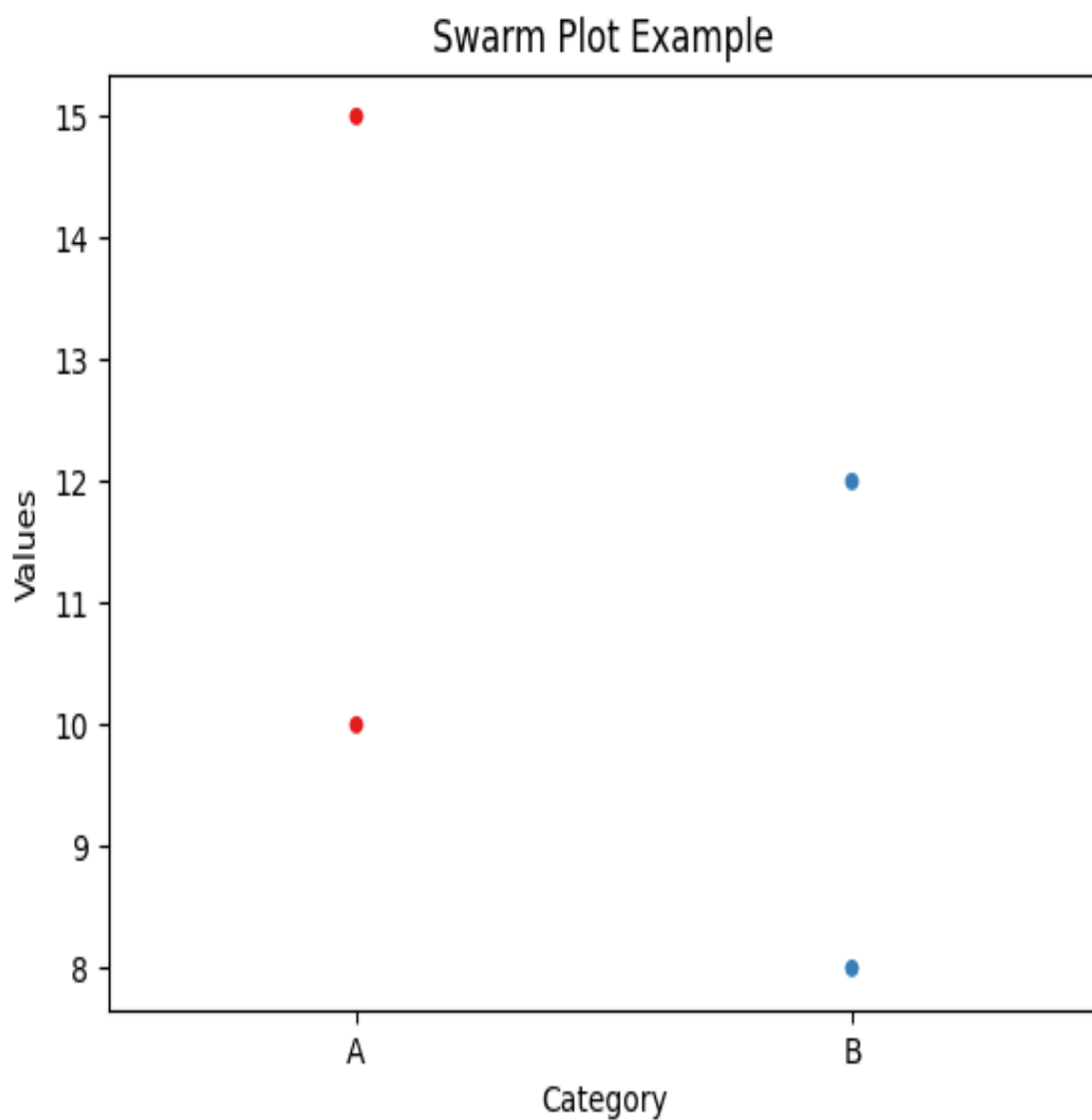
## Example 4: Violin Plot

```python
# Example data
data = {
    'Category': ['A', 'A', 'B', 'B'],
    'Values': [10, 15, 8, 12]
}
df = pd.DataFrame(data)

# Violin plot
sns.violinplot(x='Category', y='Values', data=df, palette='muted')
plt.title('Violin Plot Example')
plt.show()
```
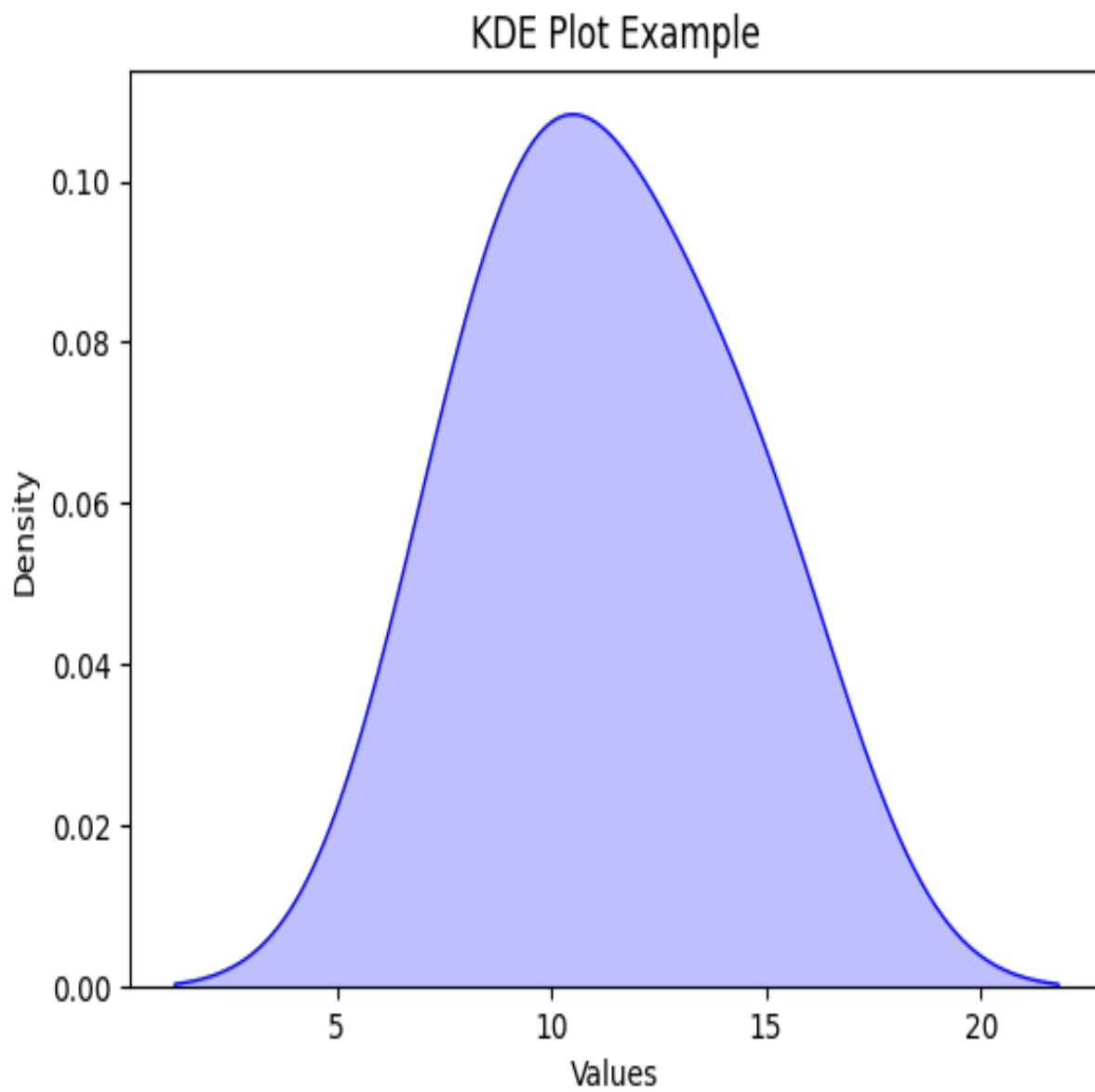
## Example 5: Swarm Plot

```python
# Example data
sns.swarmplot(x='Category', y='Values', data=df, palette='Set1')
plt.title('Swarm Plot Example')
plt.show()
```
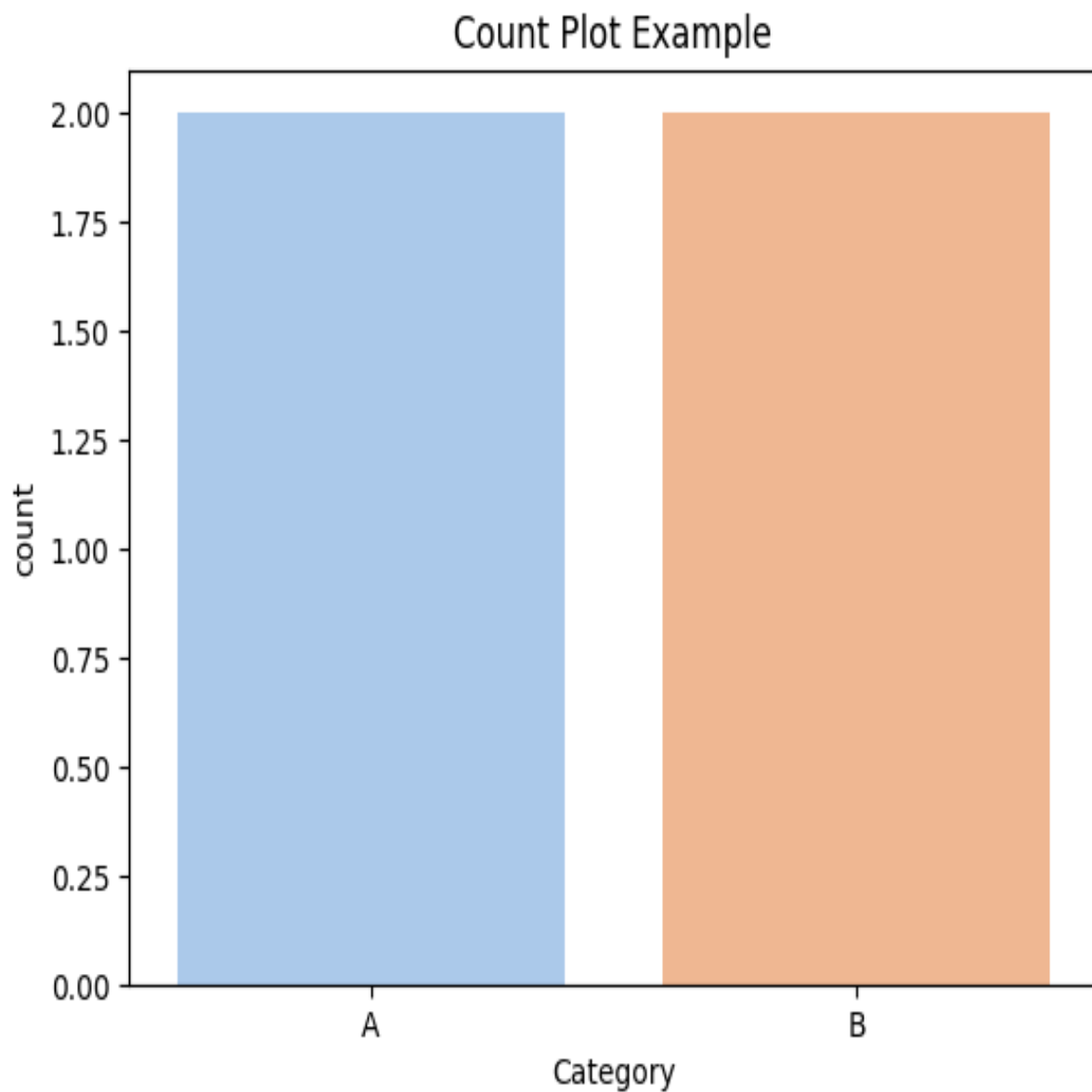


Swarm Plot Example

## Example 6: KDE Plot

```python
# KDE plot
sns.kdeplot(data=df['Values'], shade=True, color='blue')
plt.title('KDE Plot Example')
plt.show()
```



KDE Plot Example

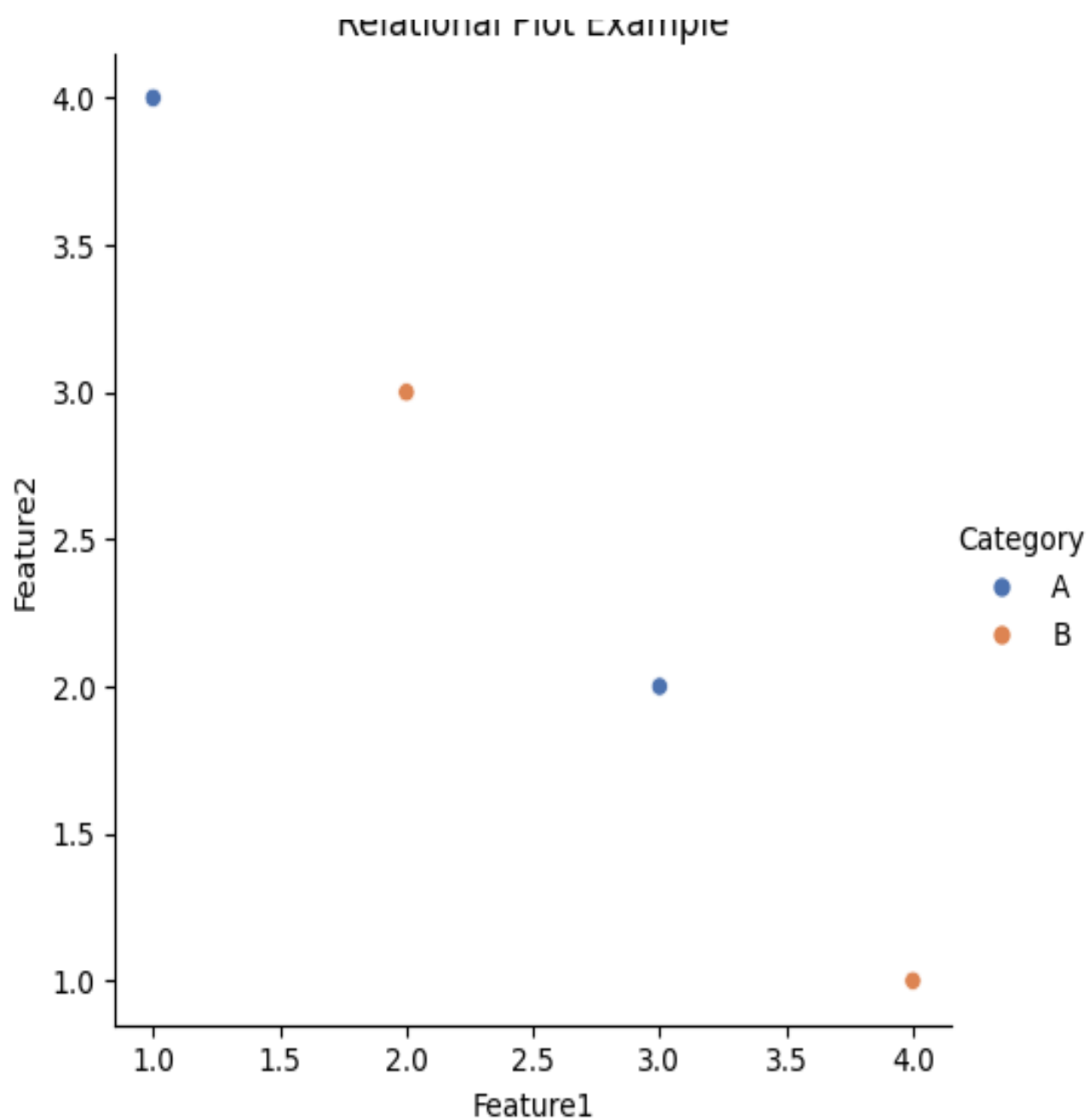## Example 7: Count Plot

```python
# Example data
sns.countplot(x='Category', data=df, palette='pastel')
plt.title('Count Plot Example')
plt.show()
```

## Example 8: Relational Plot

```
# Relational plot
sns.relplot(x='Feature1', y='Feature2', hue='Category', data=df,
kind='scatter', palette='deep')
plt.title('Relational Plot Example')
plt.show()
```

Relational Plot Example

# Chapter 12: Advanced Data Preprocessing Techniques

Effective data preprocessing is essential for accurate analysis and machine learning. In this chapter, we delve into advanced preprocessing techniques using Pandas and other Python libraries.

## Section 12.1: Handling Missing Data

### Example 1: Filling Missing Values

```python
import pandas as pd

# Example data
data = {
    'Name': ['Alice', 'Bob', 'Charlie', 'David'],
    'Age': [25, None, 30, 35],
    'Salary': [50000, 60000, None, 70000]
}
df = pd.DataFrame(data)

# Fill missing values with mean for numeric columns
df['Age'] = df['Age'].fillna(df['Age'].mean())
df['Salary'] = df['Salary'].fillna(df['Salary'].mean())
print(df)
```

**Output:**

```
      Name   Age   Salary
0    Alice  25.0  50000.0
1      Bob  30.0  60000.0
2  Charlie  30.0  60000.0
3    David  35.0  70000.0
```

### Example 2: Dropping Rows with Missing Values

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

```
print(df_dropped)
```

**Output:**

```
      Name   Age   Salary
0    Alice  25.0  50000.0
3    David  35.0  70000.0
```

# Section 12.2: Encoding Categorical Variables

## Example 1: One-Hot Encoding

```python
from sklearn.preprocessing import OneHotEncoder

# Example data
data = {'Category': ['A', 'B', 'A', 'C']}
df = pd.DataFrame(data)

# One-hot encode the category column
encoder = OneHotEncoder(sparse=False)
encoded = encoder.fit_transform(df[['Category']])
encoded_df = pd.DataFrame(encoded,
columns=encoder.get_feature_names_out(['Category']))
print(encoded_df)
```

**Output:**

```
   Category_A  Category_B  Category_C
0         1.0         0.0         0.0
1         0.0         1.0         0.0
2         1.0         0.0         0.0
3         0.0         0.0         1.0
```

## Example 2: Label Encoding

```python
from sklearn.preprocessing import LabelEncoder

# Example data
data = {'Category': ['A', 'B', 'A', 'C']}
df = pd.DataFrame(data)

# Label encode the category column
encoder = LabelEncoder()
df['Category_encoded'] = encoder.fit_transform(df['Category'])
print(df)
```

**Output:**

```
   Category  Category_encoded
0        A                 0
1        B                 1
2        A                 0
3        C                 2
```

# Section 12.3: Feature Scaling

## Example 1: Min-Max Scaling

```python
from sklearn.preprocessing import MinMaxScaler

# Example data
data = {'Feature1': [1, 2, 3, 4, 5], 'Feature2': [10, 20, 30, 40, 50]}
df = pd.DataFrame(data)

# Apply Min-Max Scaling
scaler = MinMaxScaler()
df_scaled = pd.DataFrame(scaler.fit_transform(df), columns=df.columns)
print(df_scaled)
```

**Output:**

```
   Feature1  Feature2
0      0.00      0.00
1      0.25      0.25
2      0.50      0.50
3      0.75      0.75
4      1.00      1.00
```

## Example 2: Standardization

```python
from sklearn.preprocessing import StandardScaler

# Apply Standard Scaling
scaler = StandardScaler()
df_standardized = pd.DataFrame(scaler.fit_transform(df), columns=df.columns)
print(df_standardized)
```

**Output:**

```
   Feature1  Feature2
0 -1.414214 -1.414214
1 -0.707107 -0.707107
2  0.000000  0.000000
3  0.707107  0.707107
```

4   1.414214   1.414214

# Section 12.4: Handling Outliers

## Example 1: Detecting Outliers with IQR

```python
# Calculate IQR
Q1 = df['Feature1'].quantile(0.25)
Q3 = df['Feature1'].quantile(0.75)
IQR = Q3 - Q1

# Filter outliers
df_filtered = df[(df['Feature1'] >= Q1 - 1.5 * IQR) & (df['Feature1'] <= Q3 +
1.5 * IQR)]
print(df_filtered)
```

**Output:**

```
    Feature1  Feature2
0          1        10
1          2        20
2          3        30
3          4        40
4          5        50
```

## Example 2: Clipping Outliers

```python
# Clip values outside 1st and 99th percentiles
df['Feature1'] = df['Feature1'].clip(lower=df['Feature1'].quantile(0.01),
upper=df['Feature1'].quantile(0.99))
print(df)
```

**Output:**

```
    Feature1  Feature2
0          1        10
1          2        20
2          3        30
3          4        40
4          5        50
```