

Essential Guide to Pandas

Efficient data analysis with Python.

Learn how to manipulate, analyze, and visualize data efficiently with pandas. A complete practical guide to data analysis with Python



rubió
Metabolomics

Ibon Martínez-Arranz | imartinez@labrubiocom

Data Science Manager at Rubió Metabolomics

www.rubiometabolomics.com

Ibon Martínez-Arranz, PhD in Mathematics and Statistics, holds MSc degrees in Applied Statistical Techniques and Mathematical Modeling. He has extensive expertise in statistical modeling and advanced data analysis techniques. Since 2017, he has led the Data Science area at Rubió Metabolomics, driving predictive model development and statistical computing management for metabolomics, data handling, and R&D projects. His doctoral research in Mathematics investigated and adapted genetic algorithms to improve the classification of NAFLD subtypes.

Itziar Mincholé Canals | iminchola@labrubio.com

Data Specialist at Rubió Metabolomics

www.rubiometabolomics.com

Itziar Mincholé Canals has a degree in physical sciences (1999) from the University of Zaragoza (Spain). In 1999 she began developing her professional career as a software analyst and developer, mainly working on web development projects and database programming in different sectors. In 2009 she joined OWL Metabolomics as a bioinformatician, where she has participated in several R&D projects and has developed several software tools for metabolomics. In 2016 she completed the master's degree in applied statistics with R software. After joining the Data Science department, she also supports data analysis.



Essential Guide to Pandas

Ibon Martínez-Arranz

**Life
Feels
Good**



rubió
Metabolomics



Contents

Essential Guide to Seaborn	1
Introduction	3
Data Loading	9
Read CSV File	9
Read Excel File	9
Read from SQL Database	10
Basic Data Inspection	11
Display Top Rows (df.head())	11
Display Bottom Rows (df.tail())	11
Display Data Types (df.dtypes)	12
Summary Statistics (df.describe())	12
Display Index, Columns, and Data (df.info())	12
Data Cleaning	15
Checking for Missing Values	16
Filling Missing Values	16
Dropping Missing Values	17
Renaming Columns	17
Dropping Columns	18
Data Transformation	19
Apply Function	19

Group By and Aggregate	20
Pivot Tables	21
Merge DataFrames	22
Concatenate DataFrames	23
Data Visualization Integration	25
Histogram	25
Boxplot	26
Scatter Plot	27
Line Plot	29
Bar Chart	30
Statistical Analysis	33
Correlation Matrix	33
Covariance Matrix	34
Value Counts	34
Unique Values in Column	35
Number of Unique Values	35
Indexing and Selection	37
Select Column	37
Select Multiple Columns	38
Select Rows by Position	38
Select Rows by Label	38
Conditional Selection	39
Data Formatting and Conversion	41
Convert Data Types	41
String Operations	42
Datetime Conversion	42
Setting Index	43

Advanced Data Transformation	45
Lambda Functions	45
Pivot Longer/Wider Format	46
Stack/Unstack	46
Cross Tabulations	47
Handling Time Series Data	49
Set Datetime Index	49
Resampling Data	50
Rolling Window Operations	50
File Export	53
Write to CSV	53
Write to Excel	54
Write to SQL Database	54
Advanced Data Queries	57
Query Function	57
Filtering with isin	58
Multi-Index Operations	59
Creating MultiIndex	59
Slicing on MultiIndex	60
Data Merging Techniques	61
Outer Join	61
Inner Join	62
Left Join	62
Right Join	63
Dealing with Duplicates	65
Finding Duplicates	65
Removing Duplicates	66

Custom Operations with Apply	67
Custom Apply Functions	67
Integration with Matplotlib for Custom Plots	69
Custom Plotting	69
Line Plot	69
Histogram	70
Scatter Plot	72
Bar Chart	73
Advanced Grouping and Aggregation	75
Group by Multiple Columns	75
Aggregate with Multiple Functions	76
Transform Function	76
Text Data Specific Operations	79
String Contains	79
String Split	80
Regular Expression Extraction	80
Working with JSON and XML	83
Reading JSON	83
Reading XML	84
Advanced File Handling	85
Read CSV with Specific Delimiter	85
Reading CSV with Semicolon Delimiter	85
Reading CSV with Tab Delimiter	86
Writing to JSON	87
Dealing with Missing Data	89
Interpolate Missing Values	89
Forward Fill Missing Values	90

Backward Fill Missing Values	91
Data Reshaping	93
Wide to Long Format	93
Long to Wide Format	94
Categorical Data Operations	97
Convert Column to Categorical	97
Order Categories	98
Advanced Indexing	101
Reset Index	101
Set Multiple Indexes	102
MultiIndex Slicing	102
Efficient Computations	105
Use of eval() for Efficient Operations	105
Query Method for Filtering	106
Advanced Data Merging	107
SQL-like Joins	107
Concatenating Along a Different Axis	108
Data Quality Checks	111
Assert Statement for Data Validation	111
Checking for Missing Values	111
Real-World Case Studies: Titanic Dataset	113
Description of the Data	113
Exploratory Data Analysis (EDA)	115
Data Cleaning and Preparation	120
Survival Analysis	122
Conclusions and Applications	124
Additional Resources	125



Essential Guide to Seaborn

Introduction

Welcome to our in-depth manual on Pandas, a cornerstone Python library that is indispensable in the realms of data science and analysis. Pandas provides a rich set of tools and functions that make data analysis, manipulation, and visualization both accessible and powerful.

Pandas, short for “Panel Data”, is an open-source library that offers high-level data structures and a vast array of tools for practical data analysis in Python. It has become synonymous with data wrangling, offering the DataFrame as its central data structure, which is effectively a table or a two-dimensional, size-mutable, and potentially heterogeneous tabular data structure with labeled axes (rows and columns).

To begin using Pandas, it’s typically imported alongside NumPy, another key library for numerical computations. The conventional way to import Pandas is as follows:

```
1 import pandas as pd
2 import numpy as np
```

In this manual, we will explore the multifaceted features of Pandas, covering a wide range of functionalities that cater to the needs of data analysts and scientists. Our guide will walk you through the following key areas:

1. **Data Loading:** Learn how to efficiently import data into Pandas from different sources such as CSV files, Excel sheets, and databases.
2. **Basic Data Inspection:** Understand the structure and content of your data through simple yet powerful inspection techniques.

3. **Data Cleaning:** Learn to identify and rectify inconsistencies, missing values, and anomalies in your dataset, ensuring data quality and reliability.
4. **Data Transformation:** Discover methods to reshape, aggregate, and modify data to suit your analytical needs.
5. **Data Visualization:** Integrate Pandas with visualization tools to create insightful and compelling graphical representations of your data.
6. **Statistical Analysis:** Utilize Pandas for descriptive and inferential statistics, making data-driven decisions easier and more accurate.
7. **Indexing and Selection:** Master the art of accessing and selecting data subsets efficiently for analysis.
8. **Data Formatting and Conversion:** Adapt your data into the desired format, enhancing its usability and compatibility with different analysis tools.
9. **Advanced Data Transformation:** Delve deeper into sophisticated data transformation techniques for complex data manipulation tasks.
10. **Handling Time Series Data:** Explore the handling of time-stamped data, crucial for time series analysis and forecasting.
11. **File Import/Export:** Learn how to effortlessly read from and write to various file formats, making data interchange seamless.
12. **Advanced Queries:** Employ advanced querying techniques to extract specific insights from large datasets.
13. **Multi-Index Operations:** Understand the multi-level indexing to work with high-dimensional data more effectively.
14. **Data Merging Techniques:** Explore various strategies to combine datasets, enhancing your analytical possibilities.
15. **Dealing with Duplicates:** Detect and handle duplicate records to maintain the integrity of your analysis.

16. **Custom Operations with Apply:** Harness the power of custom functions to extend Pandas' capabilities.
17. **Integration with Matplotlib for Custom Plots:** Create bespoke plots by integrating Pandas with Matplotlib, a leading plotting library.
18. **Advanced Grouping and Aggregation:** Perform complex grouping and aggregation operations for sophisticated data summaries.
19. **Text Data Specific Operations:** Manipulate and analyze textual data effectively using Pandas' string functions.
20. **Working with JSON and XML:** Handle modern data formats like JSON and XML with ease.
21. **Advanced File Handling:** Learn advanced techniques for managing file I/O operations.
22. **Dealing with Missing Data:** Develop strategies to address and impute missing values in your datasets.
23. **Data Reshaping:** Transform the structure of your data to facilitate different types of analysis.
24. **Categorical Data Operations:** Efficiently manage and analyze categorical data.
25. **Advanced Indexing:** Leverage advanced indexing techniques for more powerful data manipulation.
26. **Efficient Computations:** Optimize performance for large-scale data operations.
27. **Advanced Data Merging:** Explore sophisticated data merging and joining techniques for complex datasets.
28. **Data Quality Checks:** Implement strategies to ensure and maintain the quality of your data throughout the analysis process.

29. **Real-World Case Studies:** Apply the concepts and techniques learned throughout the manual to real-world scenarios using the Titanic dataset. This chapter demonstrates practical data analysis workflows, including data cleaning, exploratory analysis, and survival analysis, providing insights into how to utilize Pandas in practical applications to derive meaningful conclusions from complex data sets.

This manual is designed to empower you with the knowledge and skills to effectively manipulate and analyze data using Pandas, turning raw data into valuable insights. Let's begin our journey into the world of data analysis with Pandas.

Pandas, being a cornerstone in the Python data analysis landscape, has a wealth of resources and references available for those looking to delve deeper into its capabilities. Below are some key references and resources where you can find additional information, documentation, and support for working with Pandas:

1. Official Pandas Website and Documentation:

- The official website for Pandas is pandas.pydata.org. Here, you can find comprehensive documentation, including a detailed user guide, API reference, and numerous tutorials. The documentation is an invaluable resource for both beginners and experienced users, offering detailed explanations of Pandas' functionalities along with examples.

2. Pandas GitHub Repository:

- The Pandas GitHub repository, github.com/pandas-dev/pandas, is the primary source of the latest source code. It's also a hub for the development community where you can report issues, contribute to the codebase, and review upcoming features.

3. Pandas Community and Support:

- **Stack Overflow:** A large number of questions and answers can be found under the 'pandas' tag on Stack Overflow. It's a great place to seek help and contribute to community discussions.

- **Mailing List:** Pandas has an active mailing list for discussion and asking questions about usage and development.
- **Social Media:** Follow Pandas on platforms like Twitter for updates, tips, and community interactions.

4. **Scientific Python Ecosystem:**

- Pandas is a part of the larger ecosystem of scientific computing in Python, which includes libraries like NumPy, SciPy, Matplotlib, and IPython. Understanding these libraries in conjunction with Pandas can be highly beneficial.

5. **Books and Online Courses:**

- There are numerous books and online courses available that cover Pandas, often within the broader context of Python data analysis and data science. These can be excellent resources for structured learning and in-depth understanding.

6. **Community Conferences and Meetups:**

- Python and data science conferences often feature talks and workshops on Pandas. Local Python meetups can also be a good place to learn from and network with other users.

7. **Jupyter Notebooks:**

- Many online repositories and platforms host Jupyter Notebooks showcasing Pandas use cases. These interactive notebooks are excellent for learning by example and experimenting with code.

By exploring these resources, you can deepen your understanding of Pandas, stay updated with the latest developments, and connect with a vibrant community of users and contributors.

Data Loading

Efficient data loading is fundamental to any data analysis process. Pandas offers several functions to read data from different formats, making it easier to manipulate and analyze the data. In this chapter, we will explore how to read data from CSV files, Excel files, and SQL databases using Pandas.

Read CSV File

The `read_csv` function is used to load data from CSV files into a DataFrame. This function is highly customizable with numerous parameters to handle different formats and data types. Here is a basic example:

```
1 import pandas as pd
2
3 # Load data from a CSV file into a DataFrame
4 df = pd.read_csv('filename.csv')
```

This command reads data from 'filename.csv' and stores it in the DataFrame `df`. The file path can be a URL or a local file path.

Read Excel File

To read data from an Excel file, use the `read_excel` function. This function supports reading from both xls andxlsx file formats and allows you to specify the

sheet to be loaded.

```
1 # Load data from an Excel file into a DataFrame
2 df = pd.read_excel('filename.xlsx')
```

This reads the first sheet in the Excel workbook 'filename.xlsx' by default. You can specify a different sheet by using the `sheet_name` parameter.

Read from SQL Database

Pandas can also load data directly from a SQL database using the `read_sql` function. This function requires a SQL query and a connection object to the database.

```
1 import sqlalchemy
2
3 # Create a connection to a SQL database
4 engine = sqlalchemy.create_engine('sqlite:///example.db')
5 query = "SELECT * FROM my_table"
6
7 # Load data from a SQL database into a DataFrame
8 df = pd.read_sql(query, engine)
```

This example demonstrates how to connect to a SQLite database and read data from 'my_table' into a DataFrame.

Basic Data Inspection

Display Top Rows (`df.head()`)

This command, `df.head()`, displays the first five rows of the DataFrame, providing a quick glimpse of the data, including column names and some of the values.

1	A	B	C	D	E	
2	0	81	0.692744	Yes	2023-01-01	-1.082325
3	1	54	0.316586	Yes	2023-01-02	0.031455
4	2	57	0.860911	Yes	2023-01-03	-2.599667
5	3	6	0.182256	No	2023-01-04	-0.603517
6	4	82	0.210502	No	2023-01-05	-0.484947

Display Bottom Rows (`df.tail()`)

This command, `df.tail()`, shows the last five rows of the DataFrame, useful for checking the end of your dataset.

	A	B	C	D	E	
1	5	73	0.463415	No	2023-01-06	-0.442890
2	6	13	0.513276	No	2023-01-07	-0.289926
3	7	23	0.528147	Yes	2023-01-08	1.521620
4	8	87	0.138674	Yes	2023-01-09	-0.026802
5	9	39	0.005347	No	2023-01-10	-0.159331

Display Data Types (df.dtypes)

This command, `df.dtypes()`, returns the data types of each column in the DataFrame. It's helpful to understand the kind of data (integers, floats, strings, etc.) each column holds.

```
1  A          int64
2  B      float64
3  C         object
4  D  datetime64[ns]
5  E      float64
```

Summary Statistics (df.describe())

This command, `df.describe()`, provides descriptive statistics that summarize the central tendency, dispersion, and shape of a dataset's distribution, excluding NaN values. It's useful for a quick statistical overview.

```
1          A          B          E
2  count  10.000000  10.000000  10.000000
3  mean    51.500000    0.391186  -0.413633
4  std     29.963867    0.267698   1.024197
5  min      6.000000    0.005347  -2.599667
6  25%     27.000000    0.189317  -0.573874
7  50%     55.500000    0.390001  -0.366408
8  75%     79.000000    0.524429  -0.059934
9  max     87.000000    0.860911   1.521620
```

Display Index, Columns, and Data (df.info())

This command, `df.info()`, provides a concise summary of the DataFrame, including the number of non-null values in each column and the memory usage. It's

essential for initial data assessment.

```
1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 10 entries, 0 to 9
3 Data columns (total 5 columns):
4 #   Column  Non-Null Count  Dtype
5 ---  -
6 0    A      10 non-null      int64
7 1    B      10 non-null      float64
8 2    C      10 non-null      object
9 3    D      10 non-null      datetime64[ns]
10 4    E      10 non-null      float64
11 dtypes: datetime64[ns](1), float64(2), int64(1), object
    (1)
12 memory usage: 528.0 bytes
```


Data Cleaning

Let's go through the data cleaning process in a more detailed manner, step by step. We will start by creating a DataFrame that includes missing (NA or **null**) values, then apply various data cleaning operations, showing both the commands used and the resulting outputs.

First, we create a sample DataFrame that includes some missing values:

```
1 import pandas as pd
2
3 # Sample DataFrame with missing values
4 data = {
5     'old_name': [1, 2, None, 4, 5],
6     'B': [10, None, 12, None, 14],
7     'C': ['A', 'B', 'C', 'D', 'E'],
8     'D': pd.date_range(start = '2023-01-01', periods = 5,
9         freq = 'D'),
10    'E': [20, 21, 22, 23, 24]
11 }
12 df = pd.DataFrame(data)
13 df
```

	old_name	B	C	D	E
0	1.0	10.0	A	2023-01-01	20
1	2.0	NaN	B	2023-01-02	21
2	NaN	12.0	C	2023-01-03	22
3	4.0	NaN	D	2023-01-04	23
4	5.0	14.0	E	2023-01-05	24

This DataFrame contains missing values in columns 'old_name' and 'B'.

Checking for Missing Values

To find out where the missing values are located, we use:

```
1 missing_values = df.isnull().sum()
```

Result:

```
1 old_name      1
2 B             2
3 C             0
4 D             0
5 E             0
6 dtype: int64
```

Filling Missing Values

We can fill missing values with a specific value or a computed value (like the mean of the column):

```
1 filled_df = df.fillna({'old_name': 0, 'B': df['B'].mean()
                        })
```

Result:

	old_name	B	C	D	E
2	0	1.0	10.0	A	2023-01-01 20
3	1	2.0	12.0	B	2023-01-02 21
4	2	0.0	12.0	C	2023-01-03 22
5	3	4.0	12.0	D	2023-01-04 23
6	4	5.0	14.0	E	2023-01-05 24

Dropping Missing Values

Alternatively, we can drop rows with missing values:

```
1 dropped_df = df.dropna(axis = 'index')
```

Result:

	old_name	B	C	D	E	
2	0	1.0	10.0	A	2023-01-01	20
3	4	5.0	14.0	E	2023-01-05	24

We can also drop columns with missing values:

```
1 dropped_df = df.dropna(axis = 'columns')
```

Result:

	C	D	E	
2	0	A	2023-01-01	20
3	1	B	2023-01-02	21
4	2	C	2023-01-03	22
5	3	D	2023-01-04	23
6	4	E	2023-01-05	24

Renaming Columns

To rename columns for clarity or standardization:

```
1 renamed_df = df.rename(columns = {'old_name': 'A'})
```

Result:

	A	B	C	D	E	
2	0	1.0	10.0	A	2023-01-01	20
3	1	2.0	NaN	B	2023-01-02	21
4	2	NaN	12.0	C	2023-01-03	22

5	3	4.0	NaN	D	2023-01-04	23
6	4	5.0	14.0	E	2023-01-05	24

Dropping Columns

To remove unnecessary columns:

```
1 dropped_columns_df = df.drop(columns = ['E'])
```

Result:

	old_name	B	C	D
2	0	1.0	10.0	A 2023-01-01
3	1	2.0	NaN	B 2023-01-02
4	2	NaN	12.0	C 2023-01-03
5	3	4.0	NaN	D 2023-01-04
6	4	5.0	14.0	E 2023-01-05

Each of these steps demonstrates a fundamental aspect of data cleaning in Pandas, crucial for preparing your dataset for further analysis.

Data Transformation

Data transformation is a crucial step in preparing your dataset for analysis. Pandas provides powerful tools to transform, summarize, and combine data efficiently. This chapter covers key techniques such as applying functions, grouping and aggregating data, creating pivot tables, and merging or concatenating DataFrames.

Apply Function

The `apply` function allows you to apply a custom function to the DataFrame elements. This method is extremely flexible and can be applied to a single column or the entire DataFrame. Here's an example using `apply` on a single column to calculate the square of each value:

```
1 import pandas as pd
2
3 # Sample DataFrame
4 data = {'number': [1, 2, 3, 4, 5]}
5 df = pd.DataFrame(data)
6
7 # Applying a lambda function to square each value
8 df['squared'] = df['number'].apply(lambda x: x**2)
```

Result:

	number	squared
0	1	1
1	2	4

4	2	3	9
5	3	4	16
6	4	5	25

Group By and Aggregate

Grouping and aggregating data are essential for summarizing data. Here's how you can group by one column and aggregate another column using `sum`:

```
1 # Sample DataFrame
2 data = {'group': ['A', 'A', 'B', 'B', 'C'],
3         'value': [10, 15, 10, 20, 30]}
4 df = pd.DataFrame(data)
5
6 # Group by the 'group' column and sum the 'value' column
7 grouped_df = df.groupby('group').agg({'value': 'sum'})
```

Result:

	value
group	
A	25
B	30
C	30

The following Python script creates a DataFrame with data categorized by groups and two value columns. It then groups the data by the `group` column and applies different statistical aggregation functions to `value1` and `value2`. For `value1`, it calculates the mean and standard deviation. For `value2`, it computes the median and a custom measure which is a string combining the mean and standard deviation.

```
1 import pandas as pd
2 import numpy as npy
3
```

```

4 # Define a custom function that returns the mean and
  standard deviation of a series formatted as a string
5 def custom_measure(x):
6     return f"{np.mean(x):.2f} ({np.std(x):.2f})"
7
8 # Create a new DataFrame with two columns 'value'
9 data = {'group': ['A', 'A', 'B', 'B', 'C', 'C'],
10        'value1': [10, 15, 10, 20, 30, 20],
11        'value2': [5, 10, 5, 10, 15, 10]}
12 df = pd.DataFrame(data)
13
14 # Group by the 'group' column and apply different
  aggregation functions to each column
15 grouped_df = df.groupby('group').agg({
16     'value1': [('Mean', 'mean'), ('Standard Deviation', '
17     'std')], # Calculate mean and renamed standard
  deviation for value1
18     'value2': [('Median', 'median'), ('Measure',
19     'custom_measure')] # Calculate median and apply
  custom measure to value2
20 })
21
22 # Display the resulting DataFrame
23 print(grouped_df)

```

Result:

	value1	value2	
	Mean	Standard Deviation	Median
group			Measure
A	12.5	3.535534	7.5 7.50 (2.50)
B	15.0	7.071068	7.5 7.50 (2.50)
C	25.0	7.071068	12.5 12.50 (2.50)

Pivot Tables

Pivot tables are used to summarize and reorganize data in a DataFrame. Here's an example of creating a pivot table to find the mean values:


```
1 # Sample DataFrame
2 data = {'category': ['A', 'A', 'B', 'B', 'A'],
3         'value': [100, 200, 300, 400, 150]}
4 df = pd.DataFrame(data)
5
6 # Creating a pivot table
7 pivot_table = df.pivot_table(index = 'category', values =
8                               'value', aggfunc = 'mean')
```

Result:

	value
category	
A	150.0
B	350.0

Merge DataFrames

Merging DataFrames is akin to performing SQL joins. Here's an example of merging two DataFrames on a common column:

```
1 # Sample DataFrames
2 data1 = {'id': [1, 2, 3],
3         'name': ['Alice', 'Bob', 'Charlie']}
4 df1 = pd.DataFrame(data1)
5 data2 = {'id': [1, 2, 4],
6         'age': [25, 30, 35]}
7 df2 = pd.DataFrame(data2)
8
9 # Merging df1 and df2 on the 'id' column
10 merged_df = pd.merge(df1, df2, on = 'id')
```

Result:

	id	name	age
0	1	Alice	25
1	2	Bob	30

Concatenate DataFrames

Concatenating DataFrames is useful when you need to combine similar data from different sources. Here's how to concatenate two DataFrames:

```
1 # Sample DataFrames
2 data3 = {'name': ['David', 'Ella'],
3         'age': [28, 22]}
4 df3 = pd.DataFrame(data3)
5
6 # Concatenating df2 and df3
7 concatenated_df = pd.concat([df2, df3])
```

Result:

	id	age	name
2	0	1.0	25
3	1	2.0	30
4	2	4.0	35
5	0	NaN	28
6	1	NaN	22

These techniques provide a robust framework for transforming data, allowing you to prepare and analyze your datasets more effectively.

Data Visualization Integration

Visualizing data is a powerful way to understand and communicate the underlying patterns and relationships within your dataset. Pandas integrates seamlessly with Matplotlib, a comprehensive library for creating static, animated, and interactive visualizations in Python. This chapter demonstrates how to use Pandas for common data visualizations.

Histogram

Histograms are used to plot the distribution of a dataset. Here's how to create a histogram from a DataFrame column:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Sample DataFrame
5 data = {'scores': [88, 76, 90, 84, 65, 79, 93, 80]}
6 df = pd.DataFrame(data)
7
8 # Creating a histogram
9 df['scores'].hist()
10 plt.title('Distribution of Scores')
11 plt.xlabel('Scores')
12 plt.ylabel('Frequency')
13 plt.show()
```

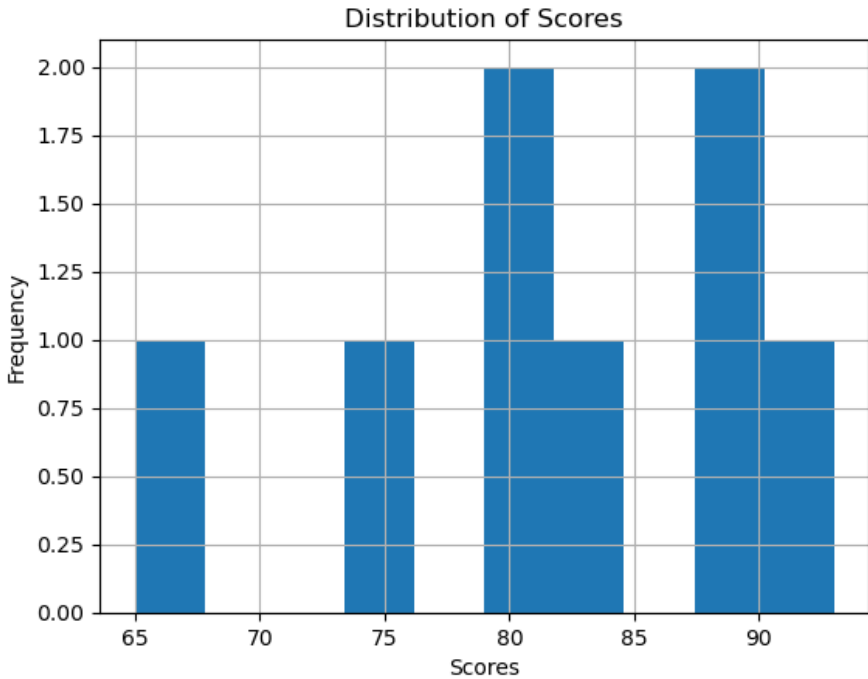


Figure 1: Image generated by the provided code.

Boxplot

Boxplots are useful for visualizing the distribution of data through their quartiles and detecting outliers. Here's how to create boxplots for multiple columns:

```
1 # Sample DataFrame
2 data = {'math_scores': [88, 76, 90, 84, 65],
3         'eng_scores': [78, 82, 88, 91, 73]}
4 df = pd.DataFrame(data)
5
6 # Creating a boxplot
```

```
7 df.boxplot(column = ['math_scores', 'eng_scores'])
8 plt.title('Score Distribution')
9 plt.ylabel('Scores')
10 plt.show()
```

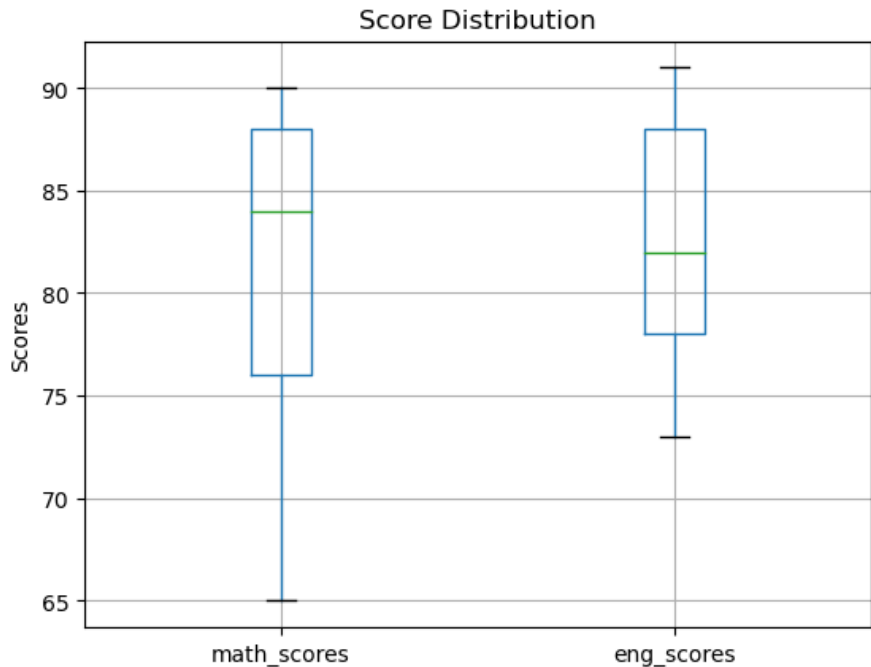


Figure 2: Image generated by the provided code.

Scatter Plot

Scatter plots are ideal for examining the relationship between two numeric variables. Here's how to create a scatter plot:

```
1 # Sample DataFrame
2 data = {'hours_studied': [10, 15, 8, 12, 6],
3         'test_score': [95, 80, 88, 90, 70]}
4 df = pd.DataFrame(data)
5
6 # Creating a scatter plot
7 df.plot.scatter(x = 'hours_studied', y = 'test_score', c
8                 = 'DarkBlue')
9 plt.title('Test Score vs Hours Studied')
10 plt.xlabel('Hours Studied')
11 plt.ylabel('Test Score')
12 plt.show()
```

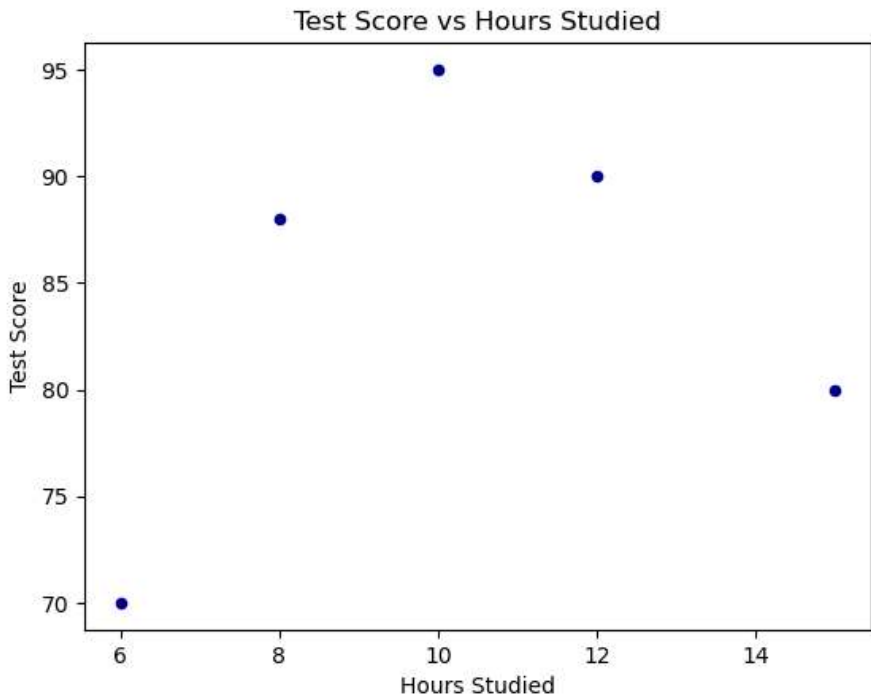


Figure 3: Image generated by the provided code.

Line Plot

Line plots are used to visualize data points connected by straight line segments. This is particularly useful in time series analysis:

```
1 # Sample DataFrame
2 data = {'year': [2010, 2011, 2012, 2013, 2014],
3         'sales': [200, 220, 250, 270, 300]}
4 df = pd.DataFrame(data)
5
6 # Creating a line plot
7 df.plot.line(x = 'year', y = 'sales', color = 'red')
8 plt.title('Yearly Sales')
9 plt.xlabel('Year')
10 plt.ylabel('Sales')
11 plt.show()
```

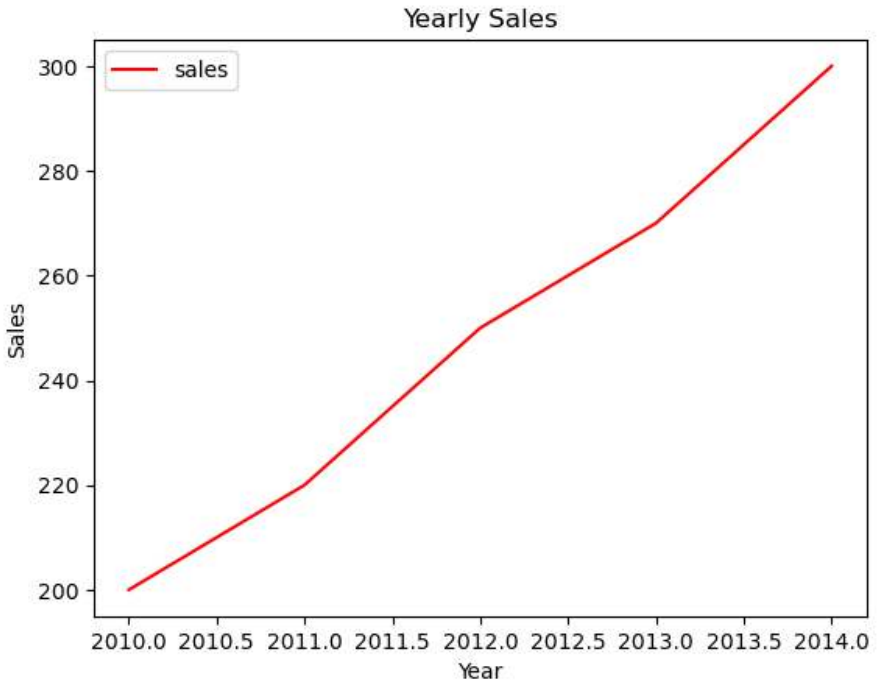



Figure 4: Image generated by the provided code.

Bar Chart

Bar charts are used to compare different groups. Here's an example of a bar chart visualizing the count of values in a column:

```
1 # Sample DataFrame
2 data = {'product': ['Apples', 'Oranges', 'Bananas', 'Apples', 'Oranges', 'Apples']}
3 df = pd.DataFrame(data)
4
5 # Creating a bar chart
```

```
6 df['product'].value_counts().plot.bar(color = 'green')
7 plt.title('Product Frequency')
8 plt.xlabel('Product')
9 plt.ylabel('Frequency')
10 plt.show()
```

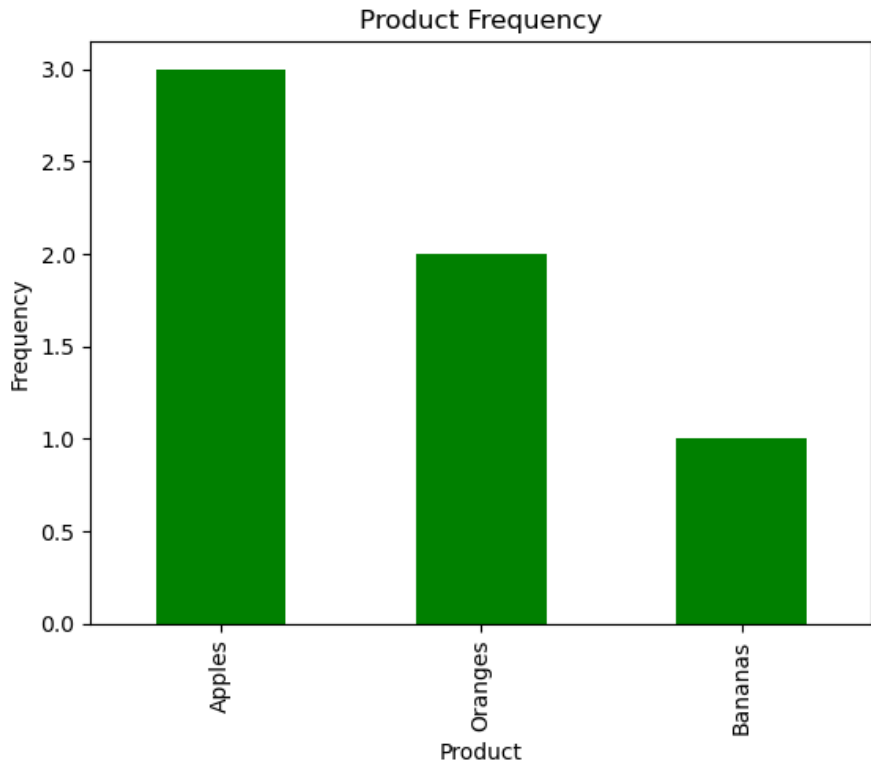


Figure 5: Image generated by the provided code.

Each of these visualization techniques provides insights into different aspects of your data, making it easier to perform comprehensive data analysis and interpretation.

Statistical Analysis

Statistical analysis is a key component of data analysis, helping to understand trends, relationships, and distributions in data. Pandas offers a range of functions for performing statistical analyses, which can be incredibly insightful when exploring your data. This chapter will cover the basics, including correlation, covariance, and various ways of summarizing data distributions.

Correlation Matrix

A correlation matrix displays the correlation coefficients between variables. Each cell in the table shows the correlation between two variables. Here's how to generate a correlation matrix:

```
1 import pandas as pd
2
3 # Sample DataFrame
4 data = {'age': [25, 30, 35, 40, 45],
5         'salary': [50000, 44000, 58000, 62000, 66000]}
6 df = pd.DataFrame(data)
7
8 # Creating a correlation matrix
9 corr_matrix = df.corr()
10 print(corr_matrix)
```

Result:

1	age	salary
---	-----	--------

```
2 age      1.000000  0.883883
3 salary   0.883883  1.000000
```

Covariance Matrix

The covariance matrix is similar to a correlation matrix but shows the covariance between variables. Here's how to generate a covariance matrix:

```
1 # Creating a covariance matrix
2 cov_matrix = df.cov()
3 print(cov_matrix)
```

Result:

```
1          age      salary
2 age      62.5      6250.0
3 salary   6250.0  80000000.0
```

Value Counts

This function is used to count the number of unique entries in a column, which can be particularly useful for categorical data:

```
1 # Sample DataFrame
2 data = {'department': ['HR', 'Finance', 'IT', 'HR', 'Finance']}
3 df = pd.DataFrame(data)
4
5 # Using value counts
6 value_counts = df['department'].value_counts()
7 print(value_counts)
```

Result:

```
1 Finance    2
2 HR         2
3 IT         1
```

Unique Values in Column

To find unique values in a column, use the `unique` function. This can help identify the diversity of entries in a column:

```
1 # Getting unique values from the column
2 unique_values = df['department'].unique()
3 print(unique_values)
```

Result:

```
1 ['HR' 'Finance' 'IT']
```

Number of Unique Values

If you need to know how many unique values are in a column, use `nunique`:

```
1 # Counting unique values
2 num_unique_values = df['department'].nunique()
3 print(num_unique_values)
```

Result:

```
1 3
```

These tools provide a fundamental insight into the statistical characteristics of your data, essential for both preliminary data exploration and advanced analyses.

Indexing and Selection

Effective data manipulation in Pandas often involves precise indexing and selection to isolate specific data segments. This chapter demonstrates several methods to select columns and rows in a DataFrame, enabling refined data analysis.

Select Column

To select a single column from a DataFrame and return it as a Series:

```
1 import pandas as pd
2
3 # Sample DataFrame
4 data = {'name': ['Alice', 'Bob', 'Charlie'],
5         'age': [25, 30, 35]}
6 df = pd.DataFrame(data)
7
8 # Selecting a single column
9 selected_column = df['name']
10 print(selected_column)
```

Result:

```
1 0    Alice
2 1     Bob
3 2   Charlie
4 Name: name, dtype: object
```


Select Multiple Columns

To select multiple columns, use a list of column names. The result is a new DataFrame:

```
1 # Selecting multiple columns
2 selected_columns = df[['name', 'age']]
3 print(selected_columns)
```

Result:

		name	age
2	0	Alice	25
3	1	Bob	30
4	2	Charlie	35

Select Rows by Position

You can select rows based on their position using `iloc`, which is primarily integer position based:

```
1 # Selecting rows by position
2 selected_rows = df.iloc[0:2]
3 print(selected_rows)
```

Result:

		name	age
2	0	Alice	25
3	1	Bob	30

Select Rows by Label

To select rows by label index, use `loc`, which uses labels in the index:

```
1 # Selecting rows by label
2 selected_rows_by_label = df.loc[0:1]
3 print(selected_rows_by_label)
```

Result:

	name	age
0	Alice	25
1	Bob	30

Conditional Selection

For conditional selection, use a condition within brackets to filter data based on column values:

```
1 # Conditional selection
2 condition_selected = df[df['age'] > 30]
3 print(condition_selected)
```

Result:

	name	age
2	Charlie	35

This selection and indexing functionality in Pandas allows for flexible and efficient data manipulations, forming the basis of many data operations you'll perform.

Data Formatting and Conversion

Data often needs to be formatted or converted to different types to meet the requirements of various analysis tasks. Pandas provides versatile capabilities for data formatting and type conversion, allowing for effective manipulation and preparation of data. This chapter covers some essential operations for data formatting and conversion.

Convert Data Types

Changing the data type of a column in a DataFrame is often necessary during data cleaning and preparation. Use `astype` to convert the data type of a column:

```
1 import pandas as pd
2
3 # Sample DataFrame
4 data = {'age': ['25', '30', '35']}
5 df = pd.DataFrame(data)
6
7 # Converting the data type of the 'age' column to integer
8 df['age'] = df['age'].astype(int)
9 print(df['age'].dtypes)
```

Result:

```
1 int64
```

String Operations

Pandas can perform vectorized string operations on Series using `.str`. This is useful for cleaning and transforming text data:

```
1 # Sample DataFrame
2 data = {'name': ['Alice', 'Bob', 'Charlie']}
3 df = pd.DataFrame(data)
4
5 # Converting all names to lowercase
6 df['name'] = df['name'].str.lower()
7 print(df)
```

Result:

```
1      name
2  0  alice
3  1   bob
4  2 charlie
```

Datetime Conversion

Converting strings or other datetime formats into a standardized `datetime64` type is essential for time series analysis. Use `pd.to_datetime` to convert a column:

```
1 # Sample DataFrame
2 data = {'date': ['2023-01-01', '2023-01-02', '2023-01-03']}
3 df = pd.DataFrame(data)
4
5 # Converting 'date' column to datetime
6 df['date'] = pd.to_datetime(df['date'])
7 print(df['date'].dtypes)
```

Result:

```
1  datetime64[ns]
```

Setting Index

Setting a specific column as the index of a DataFrame can facilitate faster searches, better alignment, and easier access to rows:

```
1  # Sample DataFrame
2  data = {'name': ['Alice', 'Bob', 'Charlie'],
3          'age': [25, 30, 35]}
4  df = pd.DataFrame(data)
5
6  # Setting 'name' as the index
7  df.set_index('name', inplace=True)
8  print(df)
```

Result:

1		age
2	name	
3	Alice	25
4	Bob	30
5	Charlie	35

These formatting and conversion techniques are crucial for preparing your dataset for detailed analysis and ensuring compatibility across different analysis and visualization tools.

Advanced Data Transformation

Advanced data transformation involves sophisticated techniques that help in re-shaping, restructuring, and summarizing complex datasets. This chapter delves into some of the more advanced functions available in Pandas that enable detailed manipulation and transformation of data.

Lambda Functions

Lambda functions provide a quick and efficient way of applying an operation across a DataFrame. Here's how you can use `apply` with a lambda function to increment every element in the DataFrame:

```
1 import pandas as pd
2
3 # Sample DataFrame
4 data = {'A': [1, 2, 3],
5         'B': [4, 5, 6]}
6 df = pd.DataFrame(data)
7
8 # Applying a lambda function to add 1 to each element
9 df = df.apply(lambda x: x + 1)
10 print(df)
```

Result:

1	A	B
2	0	2
3	2	5
4	3	6


```
3  1  3  6
4  2  4  7
```

Pivot Longer/Wider Format

The `melt` function is used to transform data from wide format to long format, which can be more suitable for analysis:

```
1 # Example of melting a DataFrame
2 data = {'Name': ['Alice', 'Bob'],
3         'Age': [25, 30],
4         'Salary': [50000, 60000]}
5 df = pd.DataFrame(data)
6
7 # Pivoting from wider to longer format
8 df_long = df.melt(id_vars = ['Name'])
9 print(df_long)
```

Result:

	Name	variable	value	
2	0	Alice	Age	25
3	1	Bob	Age	30
4	2	Alice	Salary	50000
5	3	Bob	Salary	60000

Stack/Unstack

Stacking and unstacking are powerful for reshaping a DataFrame by pivoting the columns or the index:

```
1 # Stacking and unstacking example
2 df = pd.DataFrame(data)
3
```

```
4 # Stacking
5 stacked = df.stack()
6 print(stacked)
7
8 # Unstacking
9 unstacked = stacked.unstack()
10 print(unstacked)
```

Result for stack:

```
1 0  Name      Alice
2   Age         25
3   Salary    50000
4 1  Name      Bob
5   Age         30
6   Salary    60000
7 dtype: object
```

Result for unstack:

```
1      Name Age Salary
2 0  Alice  25  50000
3 1   Bob  30  60000
```

Cross Tabulations

Cross tabulations are used to compute a simple cross-tabulation of two (or more) factors. This can be very useful in statistics and probability analysis:

```
1 # Cross-tabulation example
2 data = {'Gender': ['Female', 'Male', 'Female', 'Male'],
3         'Handedness': ['Right', 'Left', 'Right', 'Right']}
4 df = pd.DataFrame(data)
5
6 # Creating a cross tabulation
7 crosstab = pd.crosstab(df['Gender'], df['Handedness'])
```

```
8 print(crosstab)
```

Result:

```
1 Handedness  Left  Right
2 Gender
3 Female      0     2
4 Male       1     1
```

These advanced transformations enable sophisticated handling of data structures, enhancing the ability to analyze complex datasets effectively.

Handling Time Series Data

Time series data analysis is a crucial aspect of many fields such as finance, economics, and meteorology. Pandas provides robust tools for working with time series data, allowing for detailed analysis of time-stamped information. This chapter will explore how to manipulate time series data effectively using Pandas.

Set Datetime Index

Setting a datetime index is foundational in time series analysis as it facilitates easier slicing, aggregation, and resampling of data:

```
1 import pandas as pd
2
3 # Sample DataFrame with date information
4 data = {'date': ['2023-01-01', '2023-01-02', '2023-01-03',
5               '2023-01-04'],
6         'value': [100, 110, 120, 130]}
7 df = pd.DataFrame(data)
8
9 # Converting 'date' column to datetime and setting it as index
10 df['date'] = pd.to_datetime(df['date'])
11 df = df.set_index('date')
12 print(df)
```

Result:

1	value
---	-------

```
2 date
3 2023-01-01    100
4 2023-01-02    110
5 2023-01-03    120
6 2023-01-04    130
```

Resampling Data

Resampling is a powerful method for time series data aggregation or downsampling, which changes the frequency of your data:

```
1 # Resampling the data monthly and calculating the mean
2 monthly_mean = df.resample('M').mean()
3 print(monthly_mean)
```

Result:

```
1          value
2 date
3 2023-01-31  115.0
```

Rolling Window Operations

Rolling window operations are useful for smoothing or calculating moving averages, which can help in identifying trends in time series data:

```
1 # Adding more data points for a better rolling example
2 additional_data = {'date': pd.date_range('2023-01-05',
3     periods = 5, freq = 'D'),
4     'value': [140, 150, 160, 170, 180]}
5 additional_df = pd.DataFrame(additional_data)
6 df = pd.concat([df, additional_df.set_index('date')])
7 # Calculating rolling mean with a window of 5 days
```

```
8 rolling_mean = df.rolling(window = 5).mean()  
9 print(rolling_mean)
```

Result:

	date	value
3	2023-01-01	NaN
4	2023-01-02	NaN
5	2023-01-03	NaN
6	2023-01-04	NaN
7	2023-01-05	120.0
8	2023-01-06	130.0
9	2023-01-07	140.0
10	2023-01-08	150.0
11	2023-01-09	160.0

These techniques are essential for analyzing time series data efficiently, providing the tools needed to handle trends, seasonality, and other temporal structures in data.

File Export

Once data analysis is complete, it is often necessary to export data into various formats for reporting, further analysis, or sharing. Pandas provides versatile tools to export data to different file formats, including CSV, Excel, and SQL databases. This chapter will cover how to export DataFrames to these common formats.

Write to CSV

Exporting a DataFrame to a CSV file is straightforward and one of the most common methods for data sharing:

```
1 import pandas as pd
2
3 # Sample DataFrame
4 data = {'name': ['Alice', 'Bob', 'Charlie'],
5         'age': [25, 30, 35]}
6 df = pd.DataFrame(data)
7
8 # Writing the DataFrame to a CSV file
9 df.to_csv('filename.csv', index = False) # index=False
    to avoid writing row indices
```

This function will create a CSV file named `filename.csv` in the current directory without the index column.

Write to Excel

Exporting data to an Excel file can be done using the `to_excel` method, which allows for the storage of data along with formatting that can be useful for reports:

```
1 # Writing the DataFrame to an Excel file
2 df.to_excel('filename.xlsx', index = False) # index=
    False to avoid writing row indices
```

This will create an Excel file `filename.xlsx` in the current directory.

Write to SQL Database

Pandas can also export a DataFrame directly to a SQL database, which is useful for integrating analysis results into applications or storing data in a centralized database:

```
1 import sqlalchemy
2
3 # Creating a SQL connection engine
4 engine = sqlalchemy.create_engine('sqlite:///example.db')
5         # Example using SQLite
6
7 # Writing the DataFrame to a SQL database
8 df.to_sql('table_name',
9         con = engine,
10        index = False,
11        if_exists = 'replace')
```

The `to_sql` function will create a new table named `table_name` in the specified SQL database and write the DataFrame to this table. The `if_exists='replace'` parameter will replace the table if it already exists; use `if_exists='append'` to add data to an existing table instead.

These export functionalities enhance the versatility of Pandas, allowing for seam-

less transitions between different stages of data processing and sharing.

Advanced Data Queries

Performing advanced queries on a DataFrame allows for precise data filtering and extraction, which is essential for detailed analysis. This chapter explores the use of the `query` function and the `isin` method for sophisticated data querying in Pandas.

Query Function

The `query` function allows you to filter rows based on a query expression. It's a powerful way to select data dynamically:

```
1 import pandas as pd
2
3 # Sample DataFrame
4 data = {'name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'age': [25, 30, 35, 40, 45]}
6 df = pd.DataFrame(data)
7
8 # Using query to filter data
9 filtered_df = df.query('age > 30')
10 print(filtered_df)
```

Result:

	name	age
2	Charlie	35
3	David	40

4	4	Eve	45
---	---	-----	----

This query returns all rows where the `age` is greater than 30.

Filtering with `isin`

The `isin` method is useful for filtering data rows where the column value is in a predefined list of values. It's especially useful for categorical data:

```
1 # Sample DataFrame
2 data = {'name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
3         'department': ['HR', 'Finance', 'IT', 'HR', 'IT']}
4 df = pd.DataFrame(data)
5
6 # Filtering using isin
7 filtered_df = df[df['department'].isin(['HR', 'IT'])]
8 print(filtered_df)
```

Result:

	name	department
0	Alice	HR
2	Charlie	IT
3	David	HR
4	Eve	IT

This example filters rows where the `department` column contains either 'HR' or 'IT'.

These advanced querying techniques enhance the ability to perform targeted data analysis, allowing for the extraction of specific segments of data based on complex criteria.

Multi-Index Operations

Handling high-dimensional data often requires the use of multi-level indexing, or MultiIndex, which allows you to store and manipulate data with an arbitrary number of dimensions in lower-dimensional data structures like DataFrames. This chapter covers creating a MultiIndex and performing slicing operations on such structures.

Creating MultiIndex

MultiIndexing enhances data aggregation and grouping capabilities. It allows for more complex data manipulations and more sophisticated analysis:

```
1 import pandas as pd
2
3 # Sample DataFrame
4 data = {
5     'state': ['CA', 'CA', 'NY', 'NY', 'TX', 'TX'],
6     'year': [2001, 2002, 2001, 2002, 2001, 2002],
7     'population': [34.5, 35.2, 18.9, 19.7, 20.1, 20.9]
8 }
9 df = pd.DataFrame(data)
10
11 # Creating a MultiIndex DataFrame
12 df.set_index(['state', 'year'], inplace = True)
13 print(df)
```

Result:

	state	year	population
3	CA	2001	34.5
4		2002	35.2
5	NY	2001	18.9
6		2002	19.7
7	TX	2001	20.1
8		2002	20.9

Slicing on MultiIndex

Slicing a DataFrame with a MultiIndex involves specifying the ranges for each level of the index, which can be done using the `slice` function or by specifying index values directly:

```
1 # Slicing MultiIndex DataFrame
2 sliced_df = df.loc[(slice('CA', 'NY'),)]
3 print(sliced_df)
```

Result:

	state	year	population
3	CA	2001	34.5
4		2002	35.2
5	NY	2001	18.9
6		2002	19.7

This example demonstrates slicing the DataFrame to include data from states 'CA' to 'NY' for the years 2001 and 2002.

These MultiIndex operations are essential for working with complex data structures effectively, enabling more nuanced data retrieval and manipulation.

Data Merging Techniques

Merging data is a fundamental aspect of many data analysis tasks, especially when combining information from multiple sources. Pandas provides powerful functions to merge DataFrames in a manner similar to SQL joins. This chapter will cover four primary types of merges: outer, inner, left, and right joins.

Outer Join

An outer join returns all records when there is a match in either the left or right DataFrame. If there is no match, the missing side will contain `NaN`.

```
1 import pandas as pd
2
3 # Sample DataFrames
4 data1 = {'column': ['A', 'B', 'C'],
5          'values1': [1, 2, 3]}
6 df1 = pd.DataFrame(data1)
7 data2 = {'column': ['B', 'C', 'D'],
8          'values2': [4, 5, 6]}
9 df2 = pd.DataFrame(data2)
10
11 # Performing an outer join
12 outer_joined = pd.merge(df1, df2, on = 'column', how = '
    outer')
13 print(outer_joined)
```

Result:

	column	values1	values2
2	0	A	1.0
3	1	B	2.0
4	2	C	3.0
5	3	D	NaN

Inner Join

An inner join returns records that have matching values in both DataFrames.

```
1 # Performing an inner join
2 inner_joined = pd.merge(df1, df2, on = 'column', how = '
  inner')
3 print(inner_joined)
```

Result:

	column	values1	values2
2	0	B	2
3	1	C	3

Left Join

A left join returns all records from the left DataFrame, and the matched records from the right DataFrame. The result is NaN in the right side where there is no match.

```
1 # Performing a left join
2 left_joined = pd.merge(df1, df2, on = 'column', how = '
  left')
3 print(left_joined)
```

Result:

	column	values1	values2
2	0	A	1
3	1	B	2
4	2	C	3

Right Join

A right join returns all records from the right DataFrame, and the matched records from the left DataFrame. The result is `NaN` in the left side where there is no match.

```
1 # Performing a right join
2 right_joined = pd.merge(df1, df2, on = 'column', how = '
    right')
3 print(right_joined)
```

Result:

	column	values1	values2
2	0	B	2
3	1	C	3
4	2	D	<code>NaN</code>

These data merging techniques are crucial for combining data from different sources, allowing for more comprehensive analyses by creating a unified dataset from multiple disparate sources.

Dealing with Duplicates

Duplicate data can skew analysis and lead to incorrect conclusions, making it essential to identify and handle duplicates effectively. Pandas provides straightforward tools to find and remove duplicates in your datasets. This chapter will guide you through these processes.

Finding Duplicates

The `df.duplicated()` function returns a boolean series indicating whether each row is a duplicate of a row that appeared earlier in the DataFrame. Here's how to use it:

```
1 import pandas as pd
2
3 # Sample DataFrame
4 data = {'name': ['Alice', 'Bob', 'Charlie', 'Bob', 'Charlie'],
5         'age': [25, 30, 35, 30, 35]}
6 df = pd.DataFrame(data)
7
8 # Finding duplicates
9 duplicates = df.duplicated()
10 print(duplicates)
```

Result:

```
1 0    False
```

```
2 1 False
3 2 False
4 3 True
5 4 True
6 dtype: bool
```

In this output, `True` indicates that the row is a duplicate of an earlier row in the DataFrame.

Removing Duplicates

To remove the duplicate rows from the DataFrame, use the `drop_duplicates()` function. By default, this function keeps the first occurrence and removes subsequent duplicates.

```
1 # Removing duplicates
2 df_unique = df.drop_duplicates()
3 print(df_unique)
```

Result:

```
1      name  age
2  0  Alice   25
3  1   Bob   30
4  2 Charlie  35
```

This method has removed rows 3 and 4, which were duplicates of earlier rows. You can also customize this behavior with the `keep` parameter, which can be set to `'last'` to keep the last occurrence instead of the first, or `False` to remove all duplicates entirely.

These techniques are essential for ensuring data quality, enabling accurate and reliable data analysis by maintaining only unique data entries in your DataFrame.

Custom Operations with Apply

The `apply` function in Pandas is highly versatile, allowing you to execute custom functions across an entire DataFrame or along a specified axis. This flexibility makes it indispensable for performing complex operations that are not directly supported by built-in methods. This chapter will demonstrate how to use `apply` for custom operations.

Custom Apply Functions

Using `apply` with a lambda function allows you to define inline functions to apply to each row or column of a DataFrame. Here is how you can use a custom function to process data row-wise:

```
1 import pandas as pd
2
3 # Define a custom function
4 def custom_func(x, y):
5     return x * 2 + y
6
7 # Sample DataFrame
8 data = {'col1': [1, 2, 3],
9         'col2': [4, 5, 6]}
10 df = pd.DataFrame(data)
11
12 # Applying a custom function row-wise
13 df['result'] = df.apply(lambda row: custom_func(row['col1'], row['col2']), axis = 1)
```

```
14 print(df)
```

Result:

	col1	col2	result
2	0	1	4
3	1	2	5
4	2	3	6

In this example, the `custom_func` is applied to each row of the DataFrame using `apply`. The function calculates a new value based on columns 'col1' and 'col2' for each row, and the results are stored in a new column 'result'.

This method of applying custom functions is powerful for data manipulation and transformation, allowing for operations that go beyond simple arithmetic or aggregation. It's particularly useful when you need to perform operations that are specific to your data and not provided by Pandas' built-in methods.

Integration with Matplotlib for Custom Plots

Visualizing data is a key step in data analysis, providing insights that are not apparent from raw data alone. Pandas integrates smoothly with Matplotlib, a popular plotting library in Python, to offer versatile options for data visualization. This chapter will show how to create custom plots using Pandas and Matplotlib.

Custom Plotting

Pandas' plotting capabilities are built on Matplotlib, allowing for straightforward generation of various types of plots directly from DataFrame and Series objects.

Line Plot

Here's how to create a simple line plot displaying trends over a series of values:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Sample data
5 data = {'Year': [2010, 2011, 2012, 2013, 2014],
6         'Sales': [100, 150, 200, 250, 300]}
7 df = pd.DataFrame(data)
8
```



```
9 # Plotting
10 df.plot(x = 'Year', y = 'Sales', kind = 'line')
11 plt.title('Yearly Sales')
12 plt.ylabel('Sales')
13 plt.show()
```

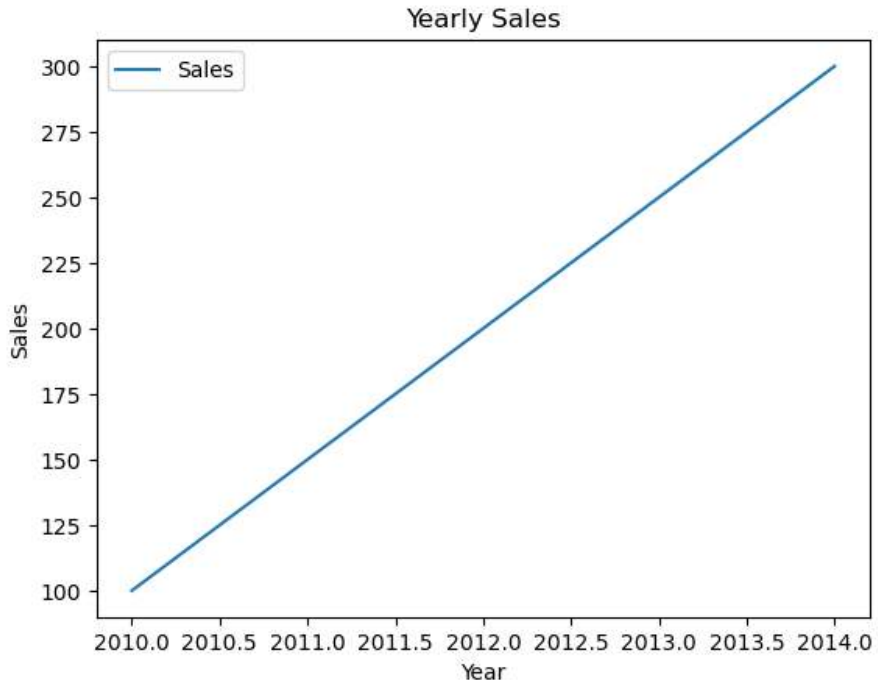


Figure 1: Image generated by the provided code.

Histogram

Histograms are great for visualizing the distribution of numerical data:

```
1 # Sample data
```

```
2 data = {'Grades': [88, 92, 80, 89, 90, 78, 84, 76, 95, 92]}
3 df = pd.DataFrame(data)
4
5 # Plotting a histogram
6 df['Grades']\
7     .plot(kind = 'hist',
8           bins = 5,
9           alpha = 0.7)
10 plt.title('Distribution of Grades')
11 plt.xlabel('Grades')
12 plt.show()
```

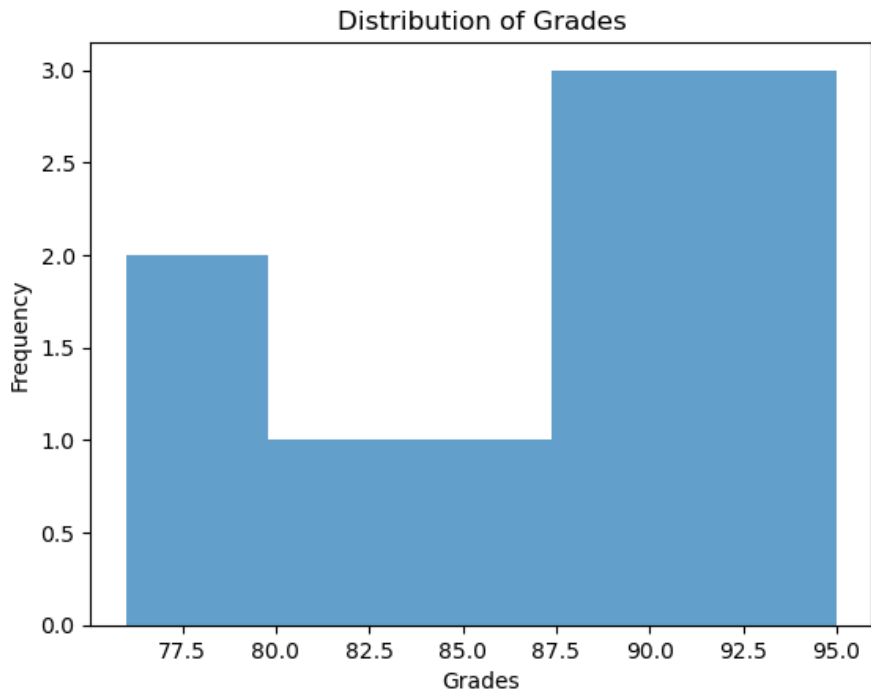


Figure 2: Image generated by the provided code.

Scatter Plot

Scatter plots are used to observe relationships between variables:

```
1 # Sample data
2 data = {'Hours': [1, 2, 3, 4, 5],
3         'Scores': [77, 78, 85, 93, 89]}
4 df = pd.DataFrame(data)
5
6 # Creating a scatter plot
7 df.plot(kind = 'scatter', x = 'Hours', y = 'Scores')
8 plt.title('Test Scores by Hours Studied')
9 plt.xlabel('Hours Studied')
10 plt.ylabel('Test Scores')
11 plt.show()
```

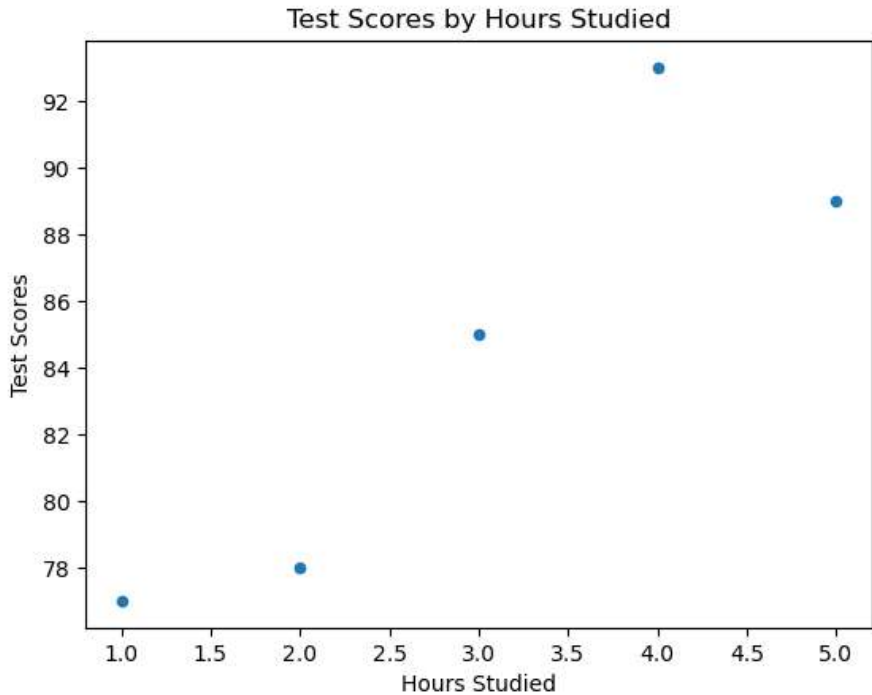


Figure 3: Image generated by the provided code.

Bar Chart

Bar charts are useful for comparing quantities corresponding to different groups:

```
1 # Sample data
2 data = {'Bars': ['A', 'B', 'C', 'D'],
3         'Values': [10, 15, 7, 10]}
4 df = pd.DataFrame(data)
5
6 # Creating a bar chart
7 df.plot(kind = 'bar',
```

```
8     x = 'Bars',  
9     y = 'Values',  
10    color = 'blue',  
11    legend = None)  
12 plt.title('Bar Chart Example')  
13 plt.ylabel('Values')  
14 plt.show()
```

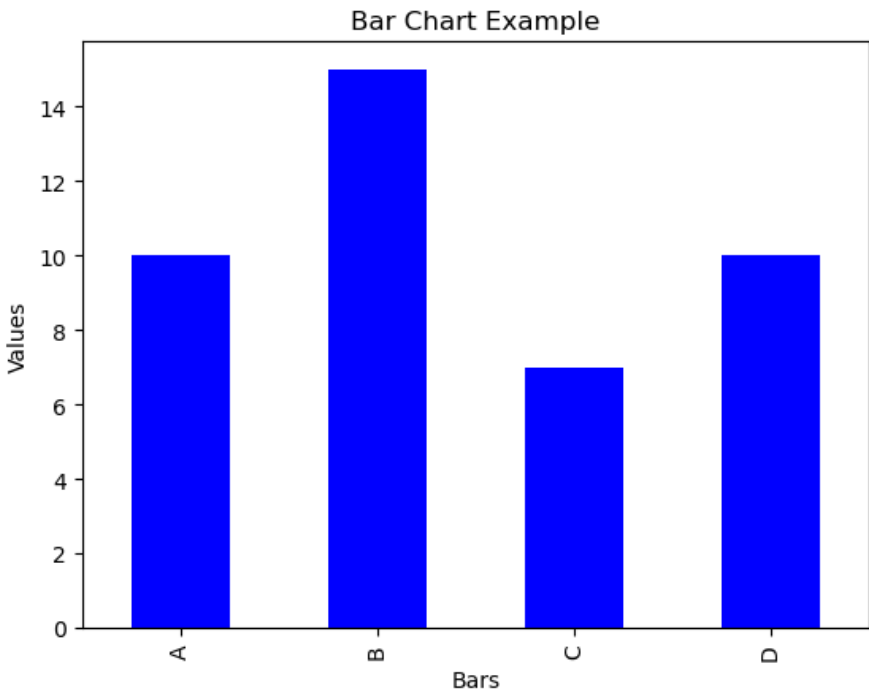


Figure 4: Image generated by the provided code.

These examples illustrate how to integrate Pandas with Matplotlib to create informative and visually appealing plots. This integration is vital for analyzing trends, distributions, relationships, and patterns in data effectively.

Advanced Grouping and Aggregation

Grouping and aggregating data are fundamental operations in data analysis, especially when dealing with large or complex datasets. Pandas offers advanced capabilities that allow for sophisticated grouping and aggregation strategies. This chapter explores some of these advanced techniques, including grouping by multiple columns, using multiple aggregation functions, and applying transformation functions.

Group by Multiple Columns

Grouping by multiple columns allows you to perform more detailed analysis. Here's how to compute the mean of groups defined by multiple columns:

```
1  import pandas as pd
2
3  # Sample DataFrame
4  data = {
5      'Department': ['Sales', 'Sales', 'HR', 'HR', 'IT', 'IT'],
6      'Team': ['A', 'B', 'A', 'B', 'A', 'B'],
7      'Revenue': [200, 210, 150, 160, 220, 230]
8  }
9  df = pd.DataFrame(data)
10
11 # Grouping by multiple columns and calculating mean
12 grouped_mean = df.groupby(['Department', 'Team']).mean()
13 print(grouped_mean)
```

Result:

		Revenue
2	Department Team	
3	HR A	150.0
4	B	160.0
5	IT A	220.0
6	B	230.0
7	Sales A	200.0
8	B	210.0

Aggregate with Multiple Functions

You can apply multiple aggregation functions at once to get a broader statistical summary:

```
1 # Applying multiple aggregation functions
2 grouped_agg = df.groupby('Department')['Revenue'].agg(['
    mean', 'sum'])
3 print(grouped_agg)
```

Result:

	Revenue	
	mean	sum
3 Department		
4 HR	155.0	310
5 IT	225.0	450
6 Sales	205.0	410

Transform Function

The `transform` function is useful for performing operations that return a DataFrame with the same index as the original. It is particularly handy for

standardizing data within groups:

```
1 # Using transform to standardize data within groups
2 df['Revenue_normalized'] = \
3     df\
4         .groupby('Department')['Revenue']\
5         .transform(lambda x: (x - x.mean()) / x.std())
6 print(df)
```

Result:

	Department	Team	Revenue	Revenue_normalized
0	Sales	A	200	-0.707107
1	Sales	B	210	0.707107
2	HR	A	150	-0.707107
3	HR	B	160	0.707107
4	IT	A	220	-0.707107
5	IT	B	230	0.707107

This example demonstrates how to normalize the 'Revenue' within each 'Department', showing deviations from the department mean in terms of standard deviations.

These advanced grouping and aggregation techniques provide powerful tools for breaking down complex data into meaningful summaries, enabling more nuanced analysis and insights.

Text Data Specific Operations

Text data often requires specific processing techniques to extract meaningful information or to reformat it for further analysis. Pandas provides a robust set of string operations that can be applied efficiently to Series and DataFrames. This chapter explores some essential operations for handling text data, including searching for substrings, splitting strings, and using regular expressions.

String Contains

The `contains` method allows you to filter rows based on whether a column's text contains a specified substring. This is useful for subsetting data based on textual content:

```
1 import pandas as pd
2
3 # Sample DataFrame
4 data = {'Description': ['Apple is sweet', 'Banana is
   yellow', 'Cherry is red']}
5 df = pd.DataFrame(data)
6
7 # Filtering rows where the Description column contains '
   sweet'
8 contains_sweet = df[df['Description'].str.contains('sweet
   ')]
9 print(contains_sweet)
```

Result:

```
1      Description
2  0  Apple is sweet
```

String Split

Splitting strings into separate components can be essential for data cleaning and preparation. The `split` method splits each string in the Series/Index by the given delimiter and optionally expands to separate columns:

```
1  # Splitting the Description column into words
2  split_description = df['Description'].str.split(' ',
3  print(split_description)
```

Result:

```
1      0      1      2
2  0  Apple  is sweet
3  1  Banana is yellow
4  2  Cherry  is   red
```

This splits the 'Description' column into separate columns for each word.

Regular Expression Extraction

Regular expressions are a powerful tool for extracting patterns from text. The `extract` method applies a regular expression pattern and extracts groups from the first match:

```
1  # Extracting the first word where it starts with a
2  # capital letter followed by lower case letters
3  extracted_words = df['Description'].str.extract(r'([A-Z][a-z]+)')
```

```
4 print(extracted_words)
```

Result:

```
1      0
2  0  Apple
3  1  Banana
4  2  Cherry
```

This regular expression extracts the first word from each description, which starts with a capital letter and is followed by lowercase letters.

These text-specific operations in Pandas simplify the process of working with textual data, allowing for efficient and powerful string manipulation and analysis.

Working with JSON and XML

In today's data-driven world, JSON (JavaScript Object Notation) and XML (eXtensible Markup Language) are two of the most common formats used for storing and transferring data on the web. Pandas provides built-in functions to easily read these formats into DataFrames, facilitating the analysis of structured data. This chapter explains how to read JSON and XML files using Pandas.

Reading JSON

JSON is a lightweight format that is easy for humans to read and write, and easy for machines to parse and generate. Pandas can directly read JSON data into a DataFrame:

```
1 import pandas as pd
2
3 # Reading JSON data
4 df = pd.read_json('filename.json')
5 print(df)
```

This method will convert a JSON file into a DataFrame. The keys of the JSON object will correspond to column names, and the values will form the data entries for the rows.

Reading XML

XML is used for representing documents with a structured markup. It is more verbose than JSON but allows for a more structured hierarchy. Pandas can read XML data into a DataFrame, similar to how it reads JSON:

```
1 # Reading XML data
2 df = pd.read_xml('filename.xml')
3 print(df)
```

This will parse an XML file and create a DataFrame. The tags of the XML file will typically define the columns, and their respective content will be the data for the rows.

These functionalities allow for seamless integration of data from web sources and other systems that utilize JSON or XML for data interchange. By leveraging Pandas' ability to work with these formats, analysts can focus more on analyzing the data rather than spending time on data preparation.

Advanced File Handling

Handling files with various configurations and formats is a common necessity in data analysis. Pandas provides extensive capabilities for reading from and writing to different file types with varying delimiters and formats. This chapter will explore reading CSV files with specific delimiters and writing DataFrames to JSON files.

Read CSV with Specific Delimiter

CSV files can come with different delimiters like commas (,), semicolons (;), or tabs (\t). Pandas allows you to specify the delimiter when reading these files, which is crucial for correctly parsing the data.

Reading CSV with Semicolon Delimiter

Suppose you have a CSV file `filename.csv` with the following content:

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

To read this CSV file into a DataFrame using Pandas, specify the semicolon as the delimiter:

```
1 import pandas as pd
```



```
2
3 # Reading a CSV file with semicolon delimiter
4 df = pd.read_csv('filename.csv', delimiter = ';')
5 print(df)
```

Result:

		Name	Age	City
2	0	Alice	30	New York
3	1	Bob	25	Los Angeles
4	2	Charlie	35	Chicago

Reading CSV with Tab Delimiter

If the CSV file uses tabs as delimiters, here's how you might see the file and read it:

File content (filename_tab.csv):

		Name	Age	City
2		Alice	30	New York
3		Bob	25	Los Angeles
4		Charlie	35	Chicago

To read this file:

```
1 # Reading a CSV file with tab delimiter
2 df_tab = pd.read_csv('filename_tab.csv', delimiter = '\t'
3 )
4 print(df_tab)
```

Result:

		Name	Age	City
2	0	Alice	30	New York
3	1	Bob	25	Los Angeles
4	2	Charlie	35	Chicago

Writing to JSON

Writing data to JSON format can be useful for web applications and APIs. Here's how to write a DataFrame to a JSON file:

```
1 # DataFrame to write to JSON
2 df.to_json('filename.json')
```

Assuming `df` contains the previous data, the JSON file `filename.json` would look like this:

```
1 {"Name":{"0":"Alice","1":"Bob","2":"Charlie"},"Age":{"0":30,"1":25,"2":35},"City":{"0":"New York","1":"Los Angeles","2":"Chicago"}}
```

This format is known as 'column-oriented' JSON. Pandas also supports other JSON orientations which can be specified using the `orient` parameter.

These advanced file handling techniques ensure that you can work with a wide range of file formats and configurations, facilitating data sharing and integration across different systems and applications.

Dealing with Missing Data

Missing data can significantly impact the results of your data analysis if not properly handled. Pandas provides several methods to deal with missing values, allowing you to either fill these gaps or make interpolations based on the existing data. This chapter explores methods like interpolation, forward filling, and backward filling.

Interpolate Missing Values

Interpolation is a method of estimating missing values by using other available data points. It is particularly useful in time series data where this can estimate the trends accurately:

```
1 import pandas as pd
2 import numpy as np
3
4 # Sample DataFrame with missing values
5 data = {'value': [1, np.nan, np.nan, 4, 5]}
6 df = pd.DataFrame(data)
7
8 # Interpolating missing values
9 df['value'] = df['value'].interpolate()
10 print(df)
```

Result:

```
1 value
```

```
2  0    1.0
3  1    2.0
4  2    3.0
5  3    4.0
6  4    5.0
```

Here, `interpolate()` linearly estimates the missing values between the existing numbers.

Forward Fill Missing Values

Forward fill (`ffill`) propagates the last observed non-null value forward until another non-null value is encountered:

```
1 # Sample DataFrame with missing values
2 data = {'value': [1, np.nan, np.nan, 4, 5]}
3 df = pd.DataFrame(data)
4
5 # Applying forward fill
6 df['value'].ffill(inplace = True)
7 print(df)
```

Result:

```
1      value
2  0    1.0
3  1    1.0
4  2    1.0
5  3    4.0
6  4    5.0
```

Backward Fill Missing Values

Backward fill (`bfill`) propagates the next observed non-null value backwards until another non-null value is met:

```
1 # Sample DataFrame with missing values
2 data = {'value': [1, np.nan, np.nan, 4, 5]}
3 df = pd.DataFrame(data)
4
5 # Applying backward fill
6 df['value'].bfill(inplace = True)
7 print(df)
```

Result:

	value
0	1.0
1	4.0
2	4.0
3	4.0
4	5.0

These methods provide you with flexible options for handling missing data based on the nature of your dataset and the specific requirements of your analysis. Correctly addressing missing data is crucial for maintaining the accuracy and reliability of your analytical results.

Data Reshaping

Data reshaping is a crucial aspect of data preparation that involves transforming data between wide format (with more columns) and long format (with more rows), depending on the needs of your analysis. This chapter demonstrates how to reshape data from wide to long formats and vice versa using Pandas.

Wide to Long Format

The `wide_to_long` function in Pandas is a powerful tool for transforming data from wide format to long format, which is often more amenable to analysis in Pandas:

```
1  import pandas as pd
2
3  # Sample DataFrame in wide format
4  data = {
5      'id': [1, 2],
6      'A_2020': [100, 200],
7      'A_2021': [150, 250],
8      'B_2020': [300, 400],
9      'B_2021': [350, 450]
10 }
11 df = pd.DataFrame(data)
12
13 # Transforming from wide to long format
14 long_df = pd.wide_to_long(df, stubnames = ['A', 'B'], sep
15                             = '_', i = 'id', j = 'year')
16 print(long_df)
```


Result:

			A	B
1				
2	id	year		
3	1	2020	100	300
4		2021	150	350
5	2	2020	200	400
6		2021	250	450

This output represents a DataFrame in long format where each row corresponds to a single year for each variable (A and B) and each id.

Long to Wide Format

Converting data from long to wide format involves creating a pivot table, which can simplify certain types of data analysis by displaying data with one variable per column and combinations of other variables per row:

```
1 # Assuming long_df is the DataFrame in long format from
  the previous example
2 # We will use a slight modification for clarity
3 long_data = {
4     'id': [1, 1, 2, 2],
5     'year': [2020, 2021, 2020, 2021],
6     'A': [100, 150, 200, 250],
7     'B': [300, 350, 400, 450]
8 }
9 long_df = pd.DataFrame(long_data)
10
11 # Transforming from long to wide format
12 wide_df = long_df.pivot(index = 'id', columns = 'year')
13 print(wide_df)
```

Result:

```
1      A      B
2  year 2020 2021 2020 2021
3  id
4  1    100  150  300  350
5  2    200  250  400  450
```

This result demonstrates a DataFrame in wide format where each `id` has associated values of A and B for each year spread across multiple columns.

Reshaping data effectively allows for easier analysis, particularly when dealing with panel data or time series that require operations across different dimensions.

Categorical Data Operations

Categorical data is common in many data sets involving categories or labels, such as survey responses, product types, or user roles. Efficient handling of such data can lead to significant performance improvements and ease of use in data manipulation and analysis. Pandas provides robust support for categorical data, including converting data types to categorical and specifying the order of categories.

Convert Column to Categorical

Converting a column to a categorical type can optimize memory usage and improve performance, especially for large datasets. Here's how to convert a column to categorical:

```
1 import pandas as pd
2
3 # Sample DataFrame
4 data = {'product': ['apple', 'banana', 'apple', 'orange',
5                    'banana', 'apple']}
5 df = pd.DataFrame(data)
6
7 # Converting 'product' column to categorical
8 df['product'] = df['product'].astype('category')
9 print(df['product'])
```

Result:

```
1 0    apple
```

```
2 1 banana
3 2 apple
4 3 orange
5 4 banana
6 5 apple
7 Name: product, dtype: category
8 Categories (3, object): ['apple', 'banana', 'orange']
```

This shows that the ‘product’ column is now of type `category` with three categories.

Order Categories

Sometimes, the natural order of categories matters (e.g., in ordinal data such as ‘low’, ‘medium’, ‘high’). Pandas allows you to set and order categories:

```
1 # Sample DataFrame with unordered categorical data
2 data = {'size': ['medium', 'small', 'large', 'small', '
    large', 'medium']}
3 df = pd.DataFrame(data)
4 df['size'] = df['size'].astype('category')
5
6 # Setting and ordering categories
7 df['size'] = df['size'].cat.set_categories(['small', '
    medium', 'large'], ordered = True, )
8 print(df['size'])
```

Result:

```
1 0 medium
2 1 small
3 2 large
4 3 small
5 4 large
6 5 medium
7 Name: size, dtype: category
8 Categories (3, object): ['small' < 'medium' < 'large']
```

This conversion and ordering process ensures that the 'size' column is not only categorical but also correctly ordered from 'small' to 'large'.

These categorical data operations in Pandas facilitate the effective handling of nominal and ordinal data, enhancing both performance and the capacity for meaningful data analysis.

Advanced Indexing

Advanced indexing techniques in Pandas enhance data manipulation capabilities, allowing for more sophisticated data retrieval and modification operations. This chapter will focus on resetting indexes, setting multiple indexes, and slicing through MultiIndexes, which are crucial for handling complex datasets effectively.

Reset Index

Resetting the index of a DataFrame can be useful when the index needs to be treated as a regular column, or when you want to revert the index back to the default integer index:

```
1 import pandas as pd
2
3 # Sample DataFrame
4 data = {'state': ['CA', 'NY', 'FL'],
5         'population': [39500000, 19500000, 21400000]}
6 df = pd.DataFrame(data)
7 df.set_index('state', inplace = True)
8
9 # Resetting the index
10 reset_df = df.reset_index(drop = True)
11 print(reset_df)
```

Result:

```
1      population
```



```
2 0 39500000
3 1 19500000
4 2 21400000
```

Using `drop=True` removes the original index and just keeps the data columns.

Set Multiple Indexes

Setting multiple columns as an index can provide powerful ways to organize and select data, especially useful in panel data or hierarchical datasets:

```
1 # Re-using previous DataFrame without resetting
2 df = pd.DataFrame(data)
3
4 # Setting multiple columns as an index
5 df.set_index(['state', 'population'], inplace = True)
6 print(df)
```

Result:

```
1 Empty DataFrame
2 Columns: []
3 Index: [(CA, 39500000), (NY, 19500000), (FL, 21400000)]
```

The DataFrame now uses a composite index made up of ‘state’ and ‘population’.

MultiIndex Slicing

Slicing data with a MultiIndex can be complex but powerful. The `xs` method (cross-section) is one of the most convenient ways to slice multi-level indexes:

```
1 # Assuming the DataFrame with a MultiIndex from the
   previous example
2 # Adding some values to demonstrate slicing
```

```
3 df['data'] = [10, 20, 30]
4
5 # Slicing with xs
6 slice_df = df.xs(key = 'CA', level = 'state')
7 print(slice_df)
```

Result:

```
1          data
2 population
3 39500000    10
```

This operation retrieves all rows associated with 'CA' from the 'state' level of the index, showing only the data for the population of California.

Advanced indexing techniques provide nuanced control over data access patterns in Pandas, enhancing data analysis and manipulation capabilities in a wide range of applications.

Efficient Computations

Efficient computation is key in handling large datasets or performing complex operations rapidly. Pandas includes features that leverage optimized code paths to speed up operations and reduce memory usage. This chapter discusses using `eval()` for arithmetic operations and the `query()` method for filtering, which are both designed to enhance performance.

Use of `eval()` for Efficient Operations

The `eval()` function in Pandas allows for the evaluation of string expressions using DataFrame columns, which can be significantly faster, especially for large DataFrames, as it avoids intermediate data copies:

```
1 import pandas as pd
2
3 # Sample DataFrame
4 data = {'col1': [1, 2, 3],
5         'col2': [4, 5, 6]}
6 df = pd.DataFrame(data)
7
8 # Using eval() to perform efficient operations
9 df['col3'] = df.eval('col1 + col2')
10 print(df)
```

Result:

1	col1	col2	col3
---	------	------	------

2	0	1	4	5
3	1	2	5	7
4	2	3	6	9

This example demonstrates how to add two columns using `eval()`, which can be faster than traditional methods for large datasets due to optimized computation.

Query Method for Filtering

The `query()` method allows you to filter DataFrame rows using an intuitive query string, which can be more readable and performant compared to traditional Boolean indexing:

```
1 # Sample DataFrame
2 data = {'col1': [10, 20, 30],
3         'col2': [20, 15, 25]}
4 df = pd.DataFrame(data)
5
6 # Using query() to filter data
7 filtered_df = df.query('col1 < col2')
8 print(filtered_df)
```

Result:

1	col1	col2	
2	0	10	20

In this example, `query()` filters the DataFrame for rows where 'col1' is less than 'col2'. This method can be especially efficient when working with large DataFrames, as it utilizes numexpr for fast evaluation of array expressions.

These methods enhance Pandas' performance, making it a powerful tool for data analysis, particularly when working with large or complex datasets. Efficient computations ensure that resources are optimally used, speeding up data processing and analysis.

Advanced Data Merging

Combining datasets is a common requirement in data analysis. Beyond basic merges, Pandas offers advanced techniques similar to SQL operations and allows concatenation along different axes. This chapter explores SQL-like joins and various concatenation methods to effectively combine multiple datasets.

SQL-like Joins

SQL-like joins in Pandas are achieved using the `merge` function. This method is extremely versatile, allowing for inner, outer, left, and right joins. Here's how to perform a left join, which includes all records from the left DataFrame and the matched records from the right DataFrame. If there is no match, the result is `NaN` on the side of the right DataFrame.

```
1 import pandas as pd
2
3 # Sample DataFrames
4 data1 = {'col': ['A', 'B', 'C'],
5          'col1': [1, 2, 3]}
6 df1 = pd.DataFrame(data1)
7 data2 = {'col': ['B', 'C', 'D'],
8          'col2': [4, 5, 6]}
9 df2 = pd.DataFrame(data2)
10
11 # Performing a left join
12 left_joined_df = pd.merge(df1, df2, how = 'left', on = '
    col')
```

```
13 print(left_joined_df)
```

Result:

	col	col1	col2
2	0	A	1
3	1	B	2
4	2	C	3

This result shows that all entries from `df1` are included, and where there are matching 'col' values in `df2`, the 'col2' values are also included.

Concatenating Along a Different Axis

Concatenation can be performed not just vertically (default axis=0), but also horizontally (axis=1). This is useful when you want to add new columns to an existing DataFrame:

```
1 # Concatenating df1 and df2 along axis 1
2 concatenated_df = pd.concat([df1, df2], axis = 1)
3 print(concatenated_df)
```

Result:

	col	col1	col	col2
2	0	A	1	B
3	1	B	2	C
4	2	C	3	D

This result demonstrates that the DataFrames are concatenated side-by-side, aligning by index. Note that because the 'col' values do not match between `df1` and `df2`, they appear disjointed, illustrating the importance of index alignment in such operations.

These advanced data merging techniques provide powerful tools for data integration, allowing for complex manipulations and combinations of datasets, much like

you would accomplish using SQL in a database environment.

Data Quality Checks

Ensuring data quality is a critical step in any data analysis process. Data often comes with issues like missing values, incorrect formats, or outliers, which can significantly impact analysis results. Pandas provides tools to perform these checks efficiently. This chapter focuses on using assertions to validate data quality.

Assert Statement for Data Validation

The `assert` statement in Python is an effective way to ensure that certain conditions are met in your data. It is used to perform sanity checks and can halt the program if the assertion fails, which is helpful in identifying data quality issues early in the data processing pipeline.

Checking for Missing Values

One common check is to ensure that there are no missing values in your DataFrame. Here's how you can use an `assert` statement to verify that there are no missing values across the entire DataFrame:

```
1 import pandas as pd
2 import numpy as np
3
4 # Sample DataFrame with possible missing values
5 data = {'col1': [1, 2, np.nan], 'col2': [4, np.nan, 6]}
6 df = pd.DataFrame(data)
```

```
7
8 # Assertion to check for missing values
9 try:
10     assert df.notnull().all().all(), "There are missing
        values in the dataframe"
11 except AssertionError as e:
12     print(e)
```

If the DataFrame contains missing values, the assertion fails, and the error message “There are missing values in the dataframe” is printed. If no missing values are present, the script continues without interruption.

This method of data validation helps in enforcing that data meets the expected quality standards before proceeding with further analysis, thus safeguarding against analysis based on faulty data.

Real-World Case Studies: Titanic Dataset

Description of the Data

This code loads the Titanic dataset directly from a publicly accessible URL into a Pandas DataFrame and prints the first few entries to get a preliminary view of the data and its structure. The `info()` function is then used to provide a concise summary of the DataFrame, detailing the non-null count and datatype for each column. This summary is invaluable for quickly identifying any missing data and understanding the data types present in each column, setting the stage for further data manipulation and analysis.

```
1 import pandas as pd
2
3 # URL of the Titanic dataset CSV from the Seaborn GitHub
  repository
4 url = "https://raw.githubusercontent.com/mwaskom/seaborn-
  data/master/titanic.csv"
5
6 # Load the dataset from the URL directly into a Pandas
  DataFrame
7 titanic = pd.read_csv(url)
8
9 # Display the first few rows of the dataframe
10 print(titanic.head())
```

```

1      survived  pclass      sex  age  sibsp  parch      fare
2  0          0        3   male  22.0     1     0    7.2500
3          S  Third
4  1          1        1  female  38.0     1     0   71.2833
5          C  First
6  2          1        3  female  26.0     0     0    7.9250
7          S  Third
8  3          1        1  female  35.0     1     0   53.1000
9          S  First
10 4          0        3   male  35.0     0     0    8.0500
11          S  Third
12
13  who  adult_male  deck  embark_town  alive  alone
14 0  man          True  NaN  Southampton  no  False
15 1  woman        False  C    Cherbourg  yes  False
16 2  woman        False  NaN  Southampton  yes  True
17 3  woman        False  C    Southampton  yes  False
18 4  man          True  NaN  Southampton  no  True

```

```

1 # Show a summary of the dataframe
2 print(titanic.info())

```

```

1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 891 entries, 0 to 890
3 Data columns (total 15 columns):
4 #   Column      Non-Null Count  Dtype
5 ---  -
6 0   survived    891 non-null    int64
7 1   pclass      891 non-null    int64
8 2   sex         891 non-null    object
9 3   age         714 non-null    float64
10 4   sibsp       891 non-null    int64
11 5   parch       891 non-null    int64
12 6   fare        891 non-null    float64
13 7   embarked    889 non-null    object
14 8   class       891 non-null    object
15 9   who         891 non-null    object
16 10  adult_male   891 non-null    bool
17 11  deck         203 non-null    object

```

```
18 12 embark_town 889 non-null object
19 13 alive      891 non-null object
20 14 alone      891 non-null bool
21 dtypes: bool(2), float64(2), int64(4), object(7)
22 memory usage: 92.4+ KB
23 None
```

Exploratory Data Analysis (EDA)

This section generates statistical summaries for numerical columns using `describe()`, which provides a quick overview of central tendencies, dispersion, and shape of the dataset's distribution. Histograms and box plots are plotted to visualize the distribution of and detect outliers in numerical data. The `value_counts()` method gives a count of unique values for categorical variables, which helps in understanding the distribution of categorical data. The `pairplot()` function from Seaborn shows pairwise relationships in the dataset, colored by the 'Survived' column to see how variables correlate with survival.

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Summary statistics for numeric columns
5 print(titanic.describe())
6
7 # Distribution of key categorical features
8 print(titanic['survived'].value_counts())
9 print(titanic['pclass'].value_counts())
10 print(titanic['sex'].value_counts())
11
12 # Histograms for numerical columns
13 titanic.hist(bins=10, figsize=(10,7))
14 plt.show()
15
16 # Box plots to check for outliers
17 titanic.boxplot(column=['age', 'fare'])
18 plt.show()
```

```

19
20 # Pairplot to visualize the relationships between
    numerical variables
21 sns.pairplot(titanic.dropna(), hue = 'survived')
22 plt.show()

```

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Summary statistics for numeric columns
5 print(titanic.describe())
6
7 # Distribution of key categorical features
8 print(titanic['survived'].value_counts())
9 print(titanic['pclass'].value_counts())
10 print(titanic['sex'].value_counts())

```

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Summary statistics for numeric columns
5 print(titanic.describe())
6
7 # Distribution of key categorical features
8 print(titanic['survived'].value_counts())
9 print(titanic['pclass'].value_counts())
10 print(titanic['sex'].value_counts())

```

	survived	pclass	age	sibsp
count	891.000000	891.000000	714.000000	891.000000
mean	0.383838	2.308642	29.699118	0.523008
std	0.486592	0.836071	14.526497	1.102743
min	0.000000	1.000000	0.420000	0.000000
25%	0.000000	2.000000	20.125000	0.000000
	0.000000	7.910400		

```
7  50%      0.000000      3.000000      28.000000      0.000000
   0.000000      14.454200
8  75%      1.000000      3.000000      38.000000      1.000000
   0.000000      31.000000
9  max      1.000000      3.000000      80.000000      8.000000
   6.000000      512.329200
10 0        549
11 1        342
12 Name: survived, dtype: int64
13 3        491
14 1        216
15 2        184
16 Name: pclass, dtype: int64
17 male      577
18 female   314
19 Name: sex, dtype: int64
```

```
1 # Histograms for numerical columns
2 titanic.hist(bins=10, figsize=(10,7))
3 plt.show()
```

```
1 # Histograms for numerical columns
2 titanic.hist(bins=10, figsize=(10,7))
3 plt.show()
```

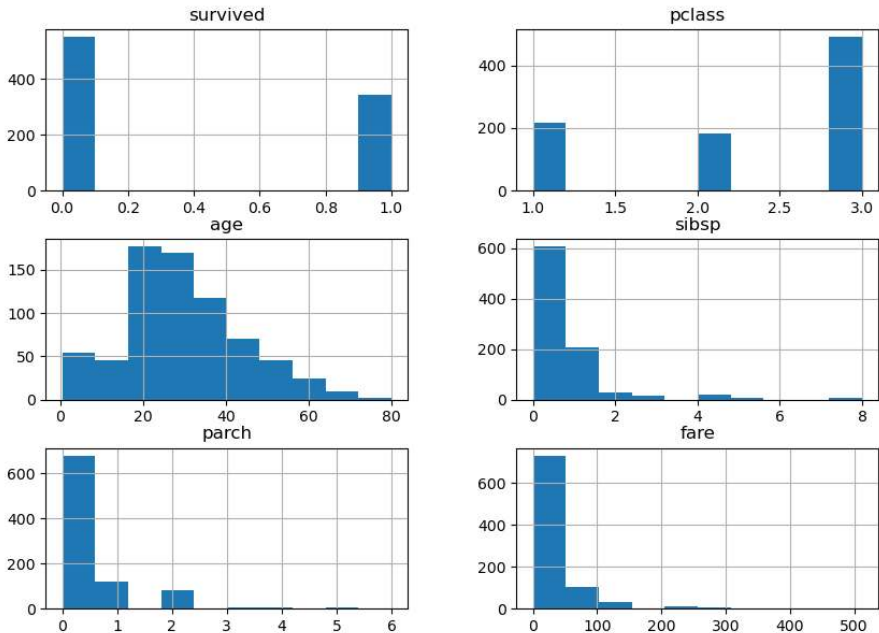



Figure 1: Image generated by the provided code.

```
1 # Box plots to check for outliers
2 titanic.boxplot(column=['age', 'fare'])
3 plt.show()
```

```
1 # Box plots to check for outliers
2 titanic.boxplot(column=['age', 'fare'])
3 plt.show()
```

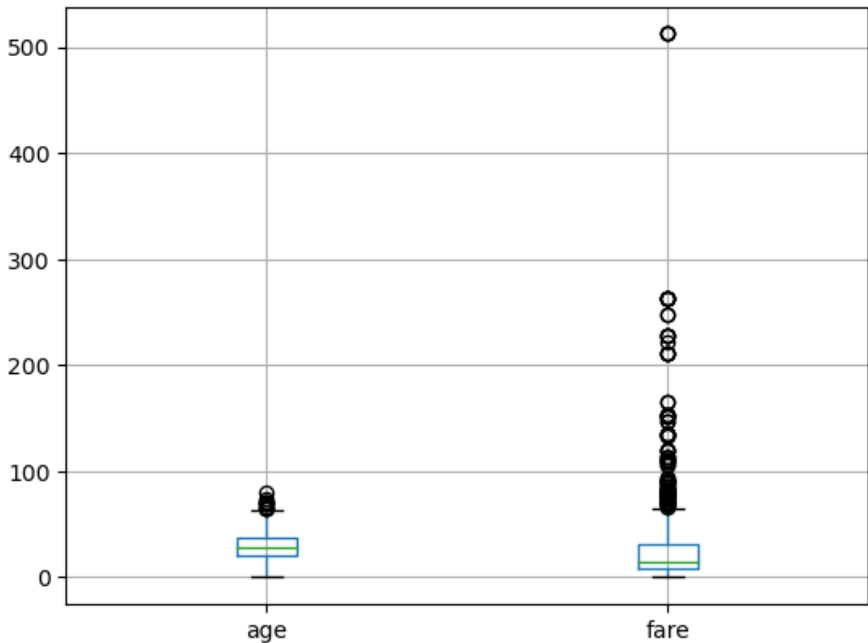


Figure 2: Image generated by the provided code.

```
1 # Pairplot to visualize the relationships between
  numerical variables
2 sns.pairplot(titanic.dropna(), hue = 'survived')
3 plt.show()
```

```
1 # Pairplot to visualize the relationships between
  numerical variables
2 sns.pairplot(titanic.dropna(), hue = 'survived')
3 plt.show()
```

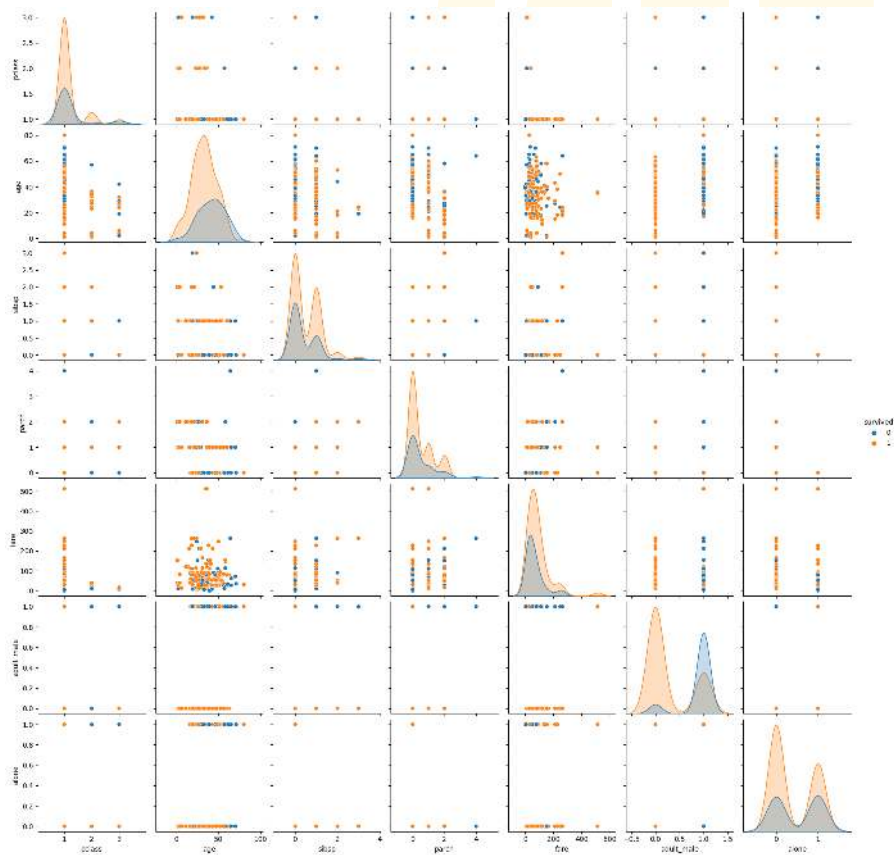


Figure 3: Image generated by the provided code.

Data Cleaning and Preparation

This code checks for missing values and handles them by filling with median values for *Age* and the mode for *Embarked*. It converts categorical data (*Sex*) into a numerical format suitable for modeling. Columns that are not necessary for the

analysis are dropped to simplify the dataset.

```
1 # Checking for missing values
2 print(titanic.isnull().sum())
3
4 # Filling missing values
5 titanic['age'].fillna(titanic['age'].median(), inplace =
   True)
6 titanic['embarked'].fillna(titanic['embarked'].mode()[0],
   inplace = True)
```

```
1 # Checking for missing values
2 print(titanic.isnull().sum())
3
4 # Filling missing values
5 titanic['age'].fillna(titanic['age'].median(), inplace =
   True)
6 titanic['embarked'].fillna(titanic['embarked'].mode()[0],
   inplace = True)
```

```
1 survived      0
2 pclass        0
3 sex           0
4 age          177
5 sibsp         0
6 parch         0
7 fare          0
8 embarked      2
9 class         0
10 who           0
11 adult_male    0
12 deck         688
13 embark_town   2
14 alive         0
15 alone         0
16 dtype: int64
```

Survival Analysis

This segment examines survival rates by *class* and *sex*. It uses `groupby()` to segment data followed by mean calculations to analyze survival rates. Results are visualized using bar plots to provide a clear visual comparison of survival rates across different groups.

```
1 # Group data by survival and class
2 survival_rate = titanic.groupby('pclass')['survived'].
    mean()
3 print(survival_rate)
```

```
1 pclass
2 1    0.629630
3 2    0.472826
4 3    0.242363
5 Name: survived, dtype: float64
```

```
1 # Survival rate by sex
2 survival_sex = titanic.groupby('sex')['survived'].mean()
3 print(survival_sex)
```

```
1 sex
2 female    0.742038
3 male      0.188908
4 Name: survived, dtype: float64
```

```
1 # Visualization of survival rates
2 sns.barplot(x = 'pclass', y = 'survived', data=titanic)
3 plt.title('Survival Rates by Class')
4 plt.show()
```

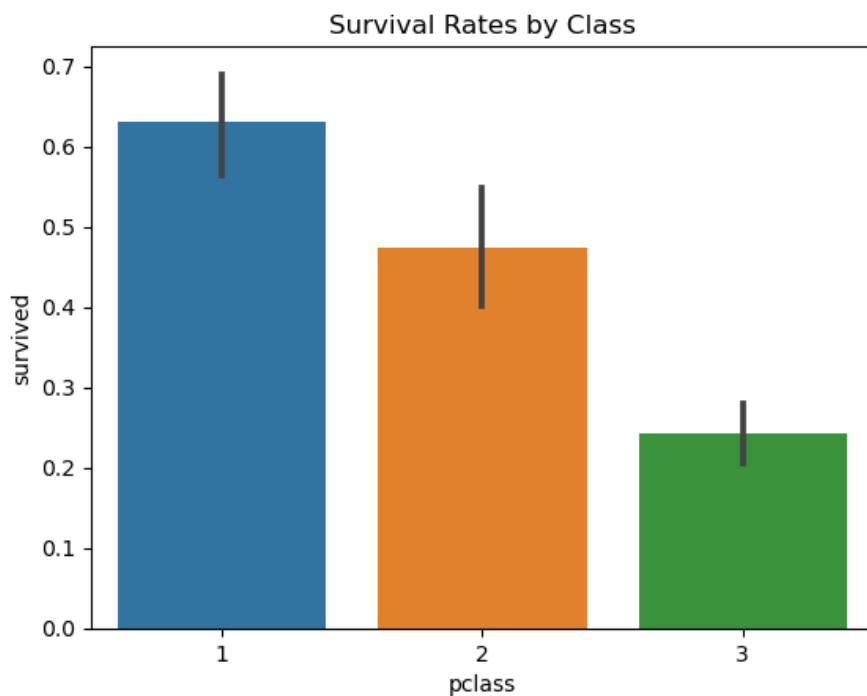


Figure 4: Image generated by the provided code.

```
1 sns.barplot(x = 'sex', y = 'survived', data=titanic)
2 plt.title('Survival Rates by Sex')
3 plt.show()
```

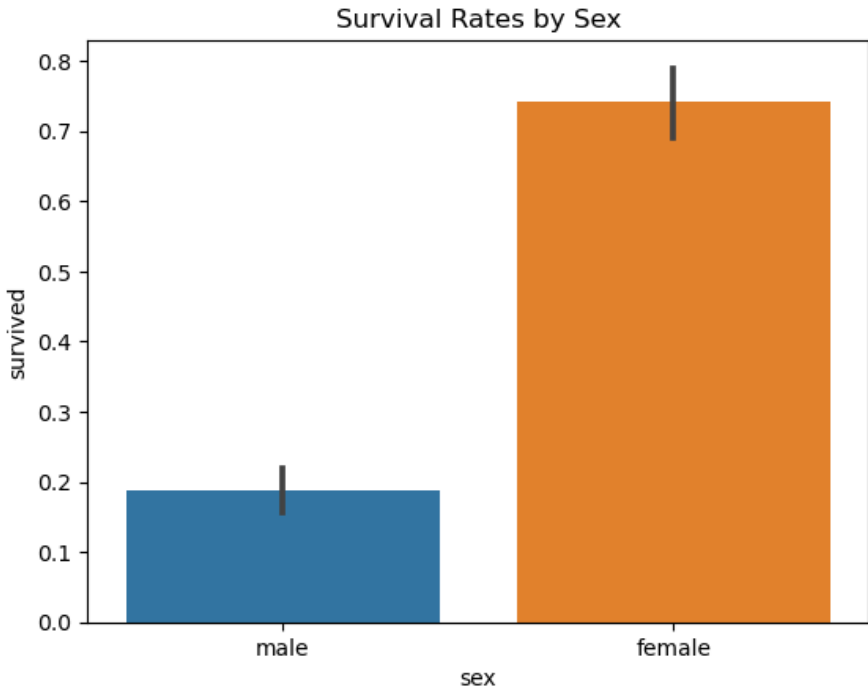


Figure 5: Image generated by the provided code.

Conclusions and Applications

The final section summarizes the key findings from the analysis, highlighting the influence of factors like *sex* and *class* on survival rates. It also discusses how the techniques applied can be used with other datasets to derive insights and support decision-making processes.

```
1 # Summary of findings
2 print("Key Findings from the Titanic Dataset:")
```

```
3 print("1. Higher survival rates were observed among
    females and upper-class passengers.")
4 print("2. Age and fare prices also appeared to influence
    survival chances.")
5
6 # Discussion on applications
7 print("These analysis techniques can be applied to other
    datasets to uncover underlying patterns and improve
    decision-making.")
```

Additional Resources

Provides additional resources for readers to explore more about [Pandas](#) and data analysis. This includes links to official documentation and the [Kaggle](#) competition page for the Titanic dataset, which offers a platform for practicing and improving data analysis skills.

This comprehensive chapter outline and code explanations give readers a thorough understanding of data analysis workflows using [Pandas](#), from data loading to cleaning, analysis, and drawing conclusions.

```
1 # This section would list URLs or references to further
  reading
2 print("For more detailed tutorials on Pandas and data
  analysis, visit:")
3 print("- The official Pandas documentation: https://
  pandas.pydata.org/pandas-docs/stable/")
4 print("- Kaggle's Titanic Competition for more
  explorations: https://www.kaggle.com/c/titanic")
```

This chapter provides a thorough walk-through using the Titanic dataset to demonstrate various data handling and analysis techniques with [Pandas](#), offering practical insights and methods that can be applied to a wide range of data analysis scenarios.

