

Pandas Data Analytics Cheat Sheet

Anis Koubaa

October 2024

Companion Notebook

For more detailed examples, visit the Google Colab notebook at: Jupyter Notebook for Pandas Data Analytics.

1. Data Loading

- **Reading Files:**

- `pd.read_csv('filename.csv')` - Load data from a CSV file.
- `pd.read_excel('filename.xlsx')` - Load data from an Excel file.
- `pd.read_sql(query, connection_object)` - Load data from a SQL table/database.
- `pd.read_json('filename.json')` - Load data from a JSON formatted string, URL or file.
- `pd.read_html('filename.html')` - Parses an html URL, string or file and extracts tables to a list of dataframes.

2. Data Cleaning

- **Handling Missing Data:**

- `df.dropna()` - Drop rows with any column having NA/null data.
- `df.fillna(value)` - Replace all NA/null data with value.

- **Drop Duplicates:**

- `df.drop_duplicates()` - Drop duplicate rows.

- **Data Type Conversion:**

- `df.astype(dtype)` - Convert the data type of a series.
- `pd.to_numeric(df['column'], errors='coerce')` - Convert data to numeric type, coerce errors.

3. Data Manipulation

- **Column Operations:**

- `df['new_column'] = df['column1'] + df['column2']` - Create a new column based on existing columns.
- `df.rename(columns={'old_name': 'new_name'})` - Rename columns.

- **Row Operations:**

- `df.sort_values('column')` - Sort the dataframe by column.
- `df.query('expression')` - Filter rows under condition.

- **Aggregation:**

- `df.groupby('column').mean()` - Group by column and return the mean of each group.
- `df.pivot_table(values='column_to_aggregate', index='row_index', columns='column_index', aggfunc=np.mean)` - Create a pivot table.

4. Data Visualization

- **Basic Plotting:**

- `df.plot(kind='line')` - Line plot of all columns.
- `df.plot(kind='bar')` - Bar plot of all columns.
- `df.plot(kind='scatter', x = 'feature1', y = 'feature2')` - Scatter plot of two features.
- `df['column'].hist()` - Histogram of a column.
- `df['column'].plot(kind='pie')` - Pie chart of a column's categories.
- `df.plot(kind='area')` - Area plot for all columns.

- **Advanced Plotting:**

- `df.boxplot(column=['col1', 'col2'])` - Box plot of one or more columns.
- `pd.plotting.scatter_matrix(df)` - Scatter matrix of DataFrame columns.
- `sns.heatmap(df.corr(), annot=True, cmap='coolwarm')` - Heatmap of the correlation matrix.
- `sns.pairplot(df)` - Pairwise plots of DataFrame's variables.
- `df['column1'].plot(kind='kde')` - Kernel Density Estimate plot for a column.

5. Complex Operations

- Merging and Joining:

- `pd.merge(df1, df2, on='key_column')` - SQL-like join operation.
- `df1.join(df2, on='key_column', how='left')` - Join matching rows from df2 to df1.

- Handling Time Series:

- `pd.to_datetime(df['date_column'])` - Convert the column to datetime.
- `df.resample('M').mean()` - Resample time-series data.

- Window Functions:

- `df.rolling(window=5).mean()` - Moving average over a window.
- `df.expanding(min_periods=1).sum()` - Cumulative sum.

Additional Tips

- Efficient Practices:

- Use `inplace=True` in functions to modify the DataFrame in place.
- Utilize vectorized operations with NumPy for performance-efficient computations.

- Miscellaneous:

- `df.describe()` - Generate descriptive statistics.
- `df.info()` - Summary of DataFrame including the index dtype and columns, non-null values, and memory usage.

1 Introduction

This document demonstrates the use of the pandas library in Python to load, clean, manipulate, and visualize data from the heart disease dataset obtained from the UCI Machine Learning Repository.

Sample Data

Below is a sample of patient data collected in a study on heart disease. We will use variations of this dataset to demonstrate merging and joining operations.

Age	Sex	CP	Trestbps	Chol	FBS	Restecg	Thalach	Exang	Oldpeak	Slope	CA	Thal	Target
63	1	1	145	233	1	2	150	0	2.3	3	0	6	0
67	1	4	160	286	0	2	108	1	1.5	2	3	3	2
67	1	4	120	229	0	2	129	1	2.6	2	2	7	1
37	1	3	130	250	0	0	187	0	3.5	3	0	3	0
41	0	2	130	204	0	2	172	0	1.4	1	0	3	0

2 Data Loading

```

1 import pandas as pd
2
3 # URL of the heart disease dataset (assuming it's named 'data.csv' locally)
4 url = 'data.csv'
5
6 # Define the column names since the dataset does not contain headers
7 column_names = ['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs', 'restecg',
8                  'thalach', 'exang', 'oldpeak', 'slope', 'ca', 'thal', 'target']
9
10 # Load the data into a DataFrame
11 heart_disease_data = pd.read_csv(url, header=None, names=column_names)
12
13 # Preview the data
14 print(heart_disease_data.head())

```

3 Data Exploration and Summary Statistics

This section demonstrates how to use pandas to explore and summarize the characteristics of a dataset, providing insights into the data's structure, content, and missing values.

3.1 Exploring Data Structure with `info()`

The `info()` method provides a concise summary of a DataFrame, including the number of non-null entries for each column, the data type of each column, and the memory usage.

```

1 # Get information about the DataFrame structure
2 heart_disease_data.info()

```

3.2 Summarizing Data Characteristics with `describe()`

The `describe()` method generates descriptive statistics that summarize the central tendency, dispersion, and shape of a dataset's distribution, excluding NaN values.

```
1 # Generate descriptive statistics
2 heart_disease_data.describe()
```

3.3 Checking for Missing Values

3.3.1 Using `isnull()` and `sum()`

To identify missing values in the dataset, `isnull()` can be used to create a Boolean series where each cell is True if data is missing. Combining `isnull()` with `sum()` provides a count of missing values per column.

```
1 # Count missing values in each column
2 heart_disease_data.isnull().sum()
```

3.4 Additional Summary Functions

- **Counting Values:** The `value_counts()` method counts the number of unique entries in a column, with the ability to normalize the results to show relative frequencies.

```
1 # Count occurrences of unique values in the 'cp' column
2 heart_disease_data['cp'].value_counts()
3
```

- **Checking for Duplicates:** The `duplicated()` function can be used to check for duplicate rows, returning a Boolean series.

```
1 # Check for duplicate rows
2 heart_disease_data.duplicated().sum()
3
```

- **Checking Unique Values:** The `nunique()` function returns the number of unique values for each column.

```
1 # Get the number of unique values in each column
2 heart_disease_data.nunique()
3
```

Each of these functions plays a crucial role in the initial stages of data analysis, helping to understand the amount of information, check for data integrity, and prepare the dataset for further analysis.

4 Data Cleaning

This section covers different data cleaning techniques including filling and removing missing values from the dataset using the pandas library.

4.1 Handling Missing Values

- Filling Missing Values:

```

1 # Fill all missing values with zero
2 heart_disease_data.fillna(0, inplace=True)
3
4 # Fill missing values with the mean of the respective column
5 heart_disease_data.fillna(heart_disease_data.mean(), inplace=True)
6
7 # Calculate the median of the 'trestbps' column
8 median_trestbps = heart_disease_data['trestbps'].median()
9
10 # Fill missing values in 'trestbps' column with its median
11 heart_disease_data['trestbps'].fillna(median_trestbps, inplace=True)
12
13 # Fill missing values forward (propagate last valid observation forward)
14 heart_disease_data.fillna(method='ffill', inplace=True)
15
16 # Fill missing values backward (use next valid observation to fill gap)
17 heart_disease_data.fillna(method='bfill', inplace=True)
18

```

- Dropping Missing Values:

```

1 # Drop rows where any cell in that row is NA
2 heart_disease_data.dropna(inplace=True)
3
4 # Drop rows where all cells in that row are NA
5 heart_disease_data.dropna(how='all', inplace=True)
6
7 # Drop columns where any cell in that column is NA
8 heart_disease_data.dropna(axis=1, inplace=True)
9
10 # Drop rows where NA values appear in specific columns (e.g., 'cp' or 'thal')
11 heart_disease_data.dropna(subset=['cp', 'thal'], inplace=True)
12

```

4.2 Handling Duplicate Entries

- Dropping Duplicate Rows:

```

1 # Drop duplicate rows
2 heart_disease_data.drop_duplicates(inplace=True)
3
4 # Drop duplicates in a subset of columns and keep the last occurrence
5 heart_disease_data.drop_duplicates(subset=['age', 'cp'], keep='last',
6 inplace=True)

```

4.3 Data Type Conversion

To ensure data integrity, it is sometimes necessary to explicitly convert the data types of certain columns. This might involve changing data types to accommodate mathematical operations, align with modeling requirements, or standardize input formats.

```

1 # Convert data type of a column to float
2 heart_disease_data['age'] = heart_disease_data['age'].astype(float)
3
4 # Convert data type of a column to integer
5 heart_disease_data['target'] = heart_disease_data['target'].astype(int)
6
7 # Use pd.to_numeric to convert the 'chol' column to numeric, coercing errors
8 heart_disease_data['chol'] = pd.to_numeric(heart_disease_data['chol'], errors='coerce')

```

3. Data Manipulation

This section details practical examples of manipulating data using pandas, including operations on columns, rows, and groups.

Column Operations

- Creating a New Column:

```

1     # Create a new column as the sum of two existing columns
2     heart_disease_data['new_column'] = heart_disease_data['age'] +
3     heart_disease_data['trestbps']

```

- Renaming Columns:

```

1     # Rename columns to more descriptive names
2     heart_disease_data.rename(columns={'cp': 'chest_pain_type', 'trestbps': 'resting_bp'}, inplace=True)
3

```

Row Operations

- Sorting Data:

```

1     # Sort the DataFrame by the 'age' column in descending order
2     heart_disease_data.sort_values('age', ascending=False, inplace=True)
3

```

- Filtering Rows:

```

1     # Filter rows where 'age' is greater than 50 and 'sex' is 1
2     filtered_data = heart_disease_data.query('age > 50 & sex == 1')
3

```

Aggregation

- Grouping and Aggregating:

```

1     # Group by 'sex' and calculate the mean of each group for 'cholesterol'
2     average_cholesterol_per_sex = heart_disease_data.groupby('sex')['chol'].mean()
3

```

- **Creating a Pivot Table:**

```

1 # Create a pivot table that calculates the average 'thalach' for each 'sex'
2     and 'cp' combination
3 pivot_table = heart_disease_data.pivot_table(values='thalach', index='sex',
4 columns='cp', aggfunc='mean')
5

```

Merging and Joining

- **SQL-like Join Operation:**

```

1 import pandas as pd
2
3 # Assume df1 and df2 are derived from the main dataset with modifications
4 df1 = data[['age', 'sex', 'cp', 'target']].copy()
5 df2 = data[['cp', 'thalach', 'exang']].drop_duplicates().copy()
6
7 # Merge two dataframes on a 'cp' (chest pain type) column
8 merged_data = pd.merge(df1, df2, on='cp')
9

```

Description: This function merges two DataFrames based on the 'cp' column, akin to a SQL join operation. This can merge detailed patient data with exercise-related outcomes, assuming 'cp' acts as a unique key for these exercises.

- **Join Matching Rows from Another DataFrame:**

```

1 # Join df2 to df1 based on the 'cp' column which is set as an index in df2
2 df1.set_index('cp', inplace=True)
3 df2.set_index('cp', inplace=True)
4 joined_data = df1.join(df2, how='left')
5

```

Description: This left join appends columns from df2 to df1 based on matching 'cp' values. It's useful for enriching patient records with corresponding exercise data without duplicating records.

4. Data Visualization

This section demonstrates various plotting techniques using pandas and seaborn to explore and represent data visually.

Basic Plotting

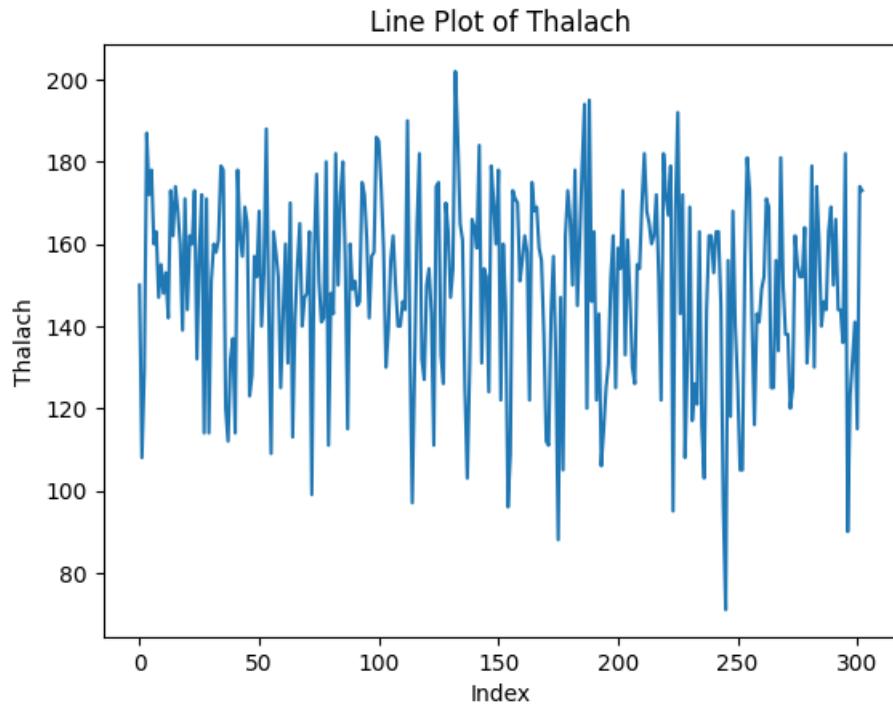
- **Line Plot:**

```

1 # Plot 'thalach' (maximum heart rate achieved) over the index
2 import matplotlib.pyplot as plt
3 heart_disease_data['thalach'].plot(kind='line')
4 plt.xlabel('Index')
5 plt.ylabel('Thalach')
6 plt.show()
7

```

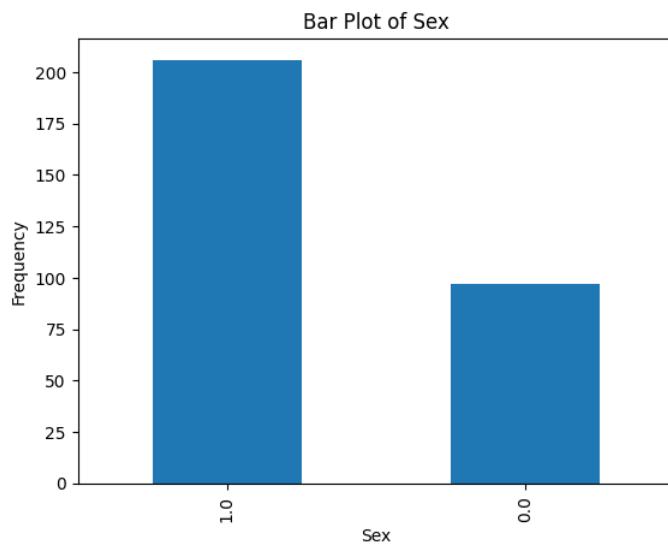
Description: A line plot is useful for showing trends over a period of time or ordinal categories. Expect to see how 'thalach' varies across the index, potentially indicating trends or patterns in heart rate data across observations.



- Bar Plot:

```
1 # Bar plot of 'sex' column to count occurrences of each category
2 data['sex'].value_counts().plot(kind='bar', title='Bar Plot of Gender')
3 plt.xlabel('Gender')
4 plt.ylabel('Frequency')
5 plt.show()
```

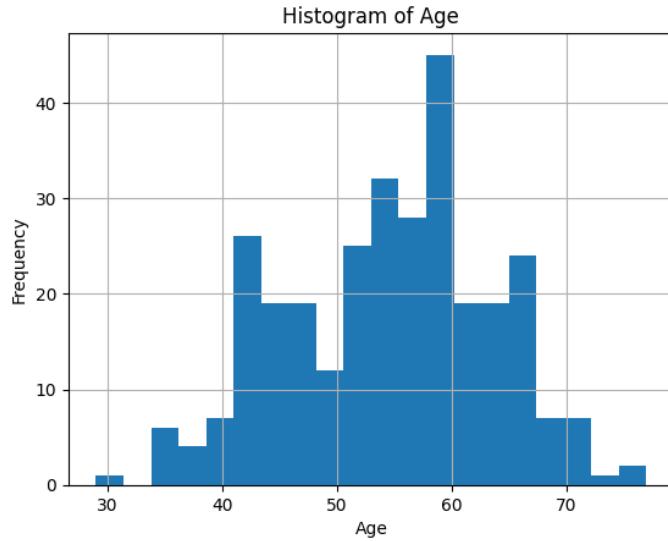
Description: Bar plots help in comparing quantities corresponding to different groups. Expect to see the distribution of data between different genders, highlighting potential imbalances or differences in sample sizes.



- **Histogram:**

```
1 # Histogram of 'age' to see the distribution of ages
2 heart_disease_data['age'].hist(bins=20)
3 plt.xlabel('Age')
4 plt.ylabel('Frequency')
5 plt.show()
```

Description: Histograms are great for visualizing the distribution of numerical data. Expect to see the frequency distribution of ages, which helps in understanding the age composition of the subjects.



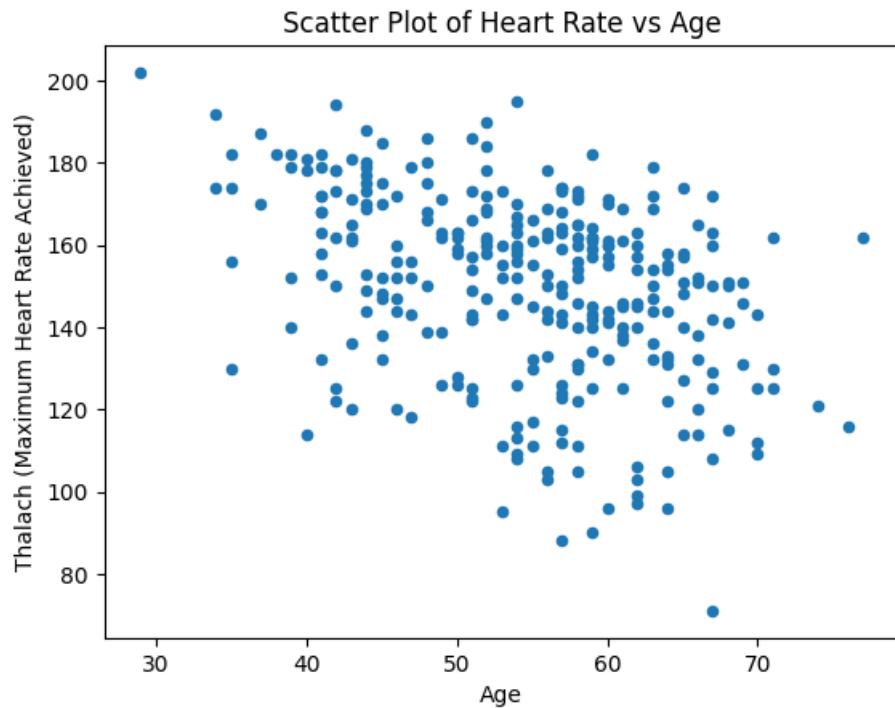
- Scatter Plot:

```

1 # Plot 'thalach' vs 'age' (maximum heart rate achieved vs age)
2 import matplotlib.pyplot as plt
3 heart_disease_data.plot(kind='scatter', x='age', y='thalach')
4 plt.xlabel('Age')
5 plt.ylabel('Thalach (Maximum Heart Rate Achieved)')
6 plt.title('Scatter Plot of Heart Rate vs Age')
7 plt.show()
8

```

Description: A scatter plot is ideal for showing the relationship between two numerical variables. This plot of 'thalach' versus 'age' helps to visualize how maximum heart rate achieved varies with age, which can be useful for identifying patterns or correlations between these variables.



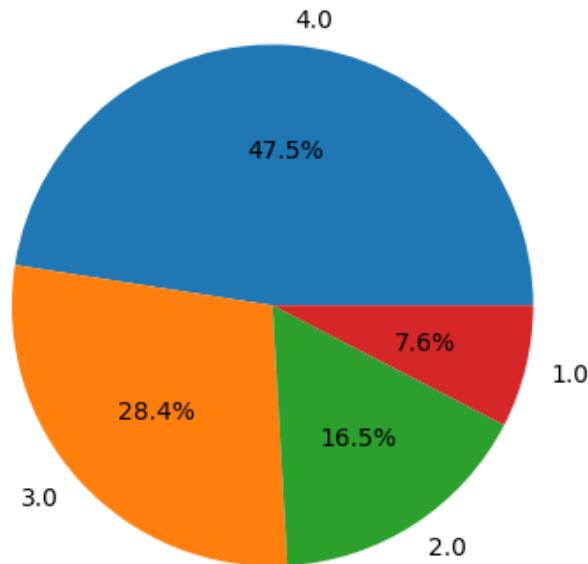
- Pie Chart:

```

1 # Pie chart of 'cp' (chest pain type)
2 heart_disease_data['cp'].value_counts().plot(kind='pie')
3 plt.ylabel('') # Remove the y-label as it's unnecessary for pie charts
4 plt.show()
5

```

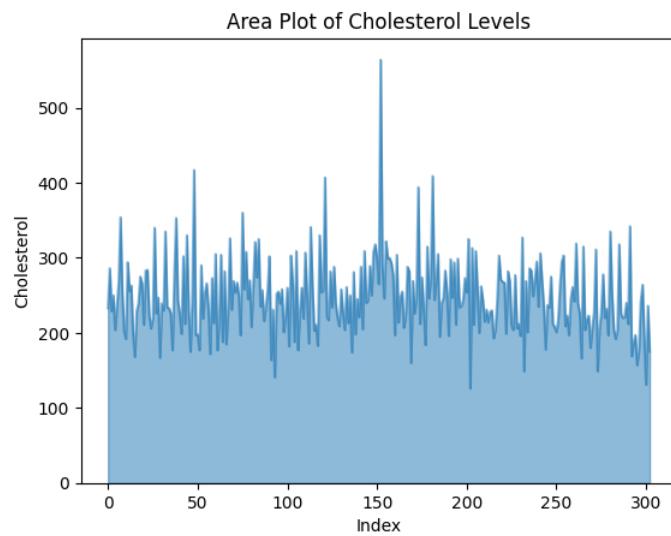
Description: Pie charts are effective for showing relative proportions or percentages of categories. Expect to see how different types of chest pain contribute to the overall dataset.

Pie Chart of Chest Pain Type

- **Area Plot:**

```
1 # Area plot for 'chol' (cholesterol levels) over the index
2 heart_disease_data['chol'].plot(kind='area', alpha=0.5)
3 plt.xlabel('Index')
4 plt.ylabel('Cholesterol')
5 plt.show()
6
```

Description: Area plots are useful for tracking changes over time among related attributes. Expect to see the area under the curve filled, indicating the magnitude of cholesterol levels through the dataset's range.



Advanced Plotting

- Box Plot:

```

1 # Box plot for 'age' and 'trestbps' (resting blood pressure)
2 heart_disease_data.boxplot(column=['age', 'trestbps'])
3 plt.title('Box Plot of Age and Resting Blood Pressure')
4 plt.ylabel('Values')
5 plt.show()
6

```

Description: Box plots are useful for identifying outliers and understanding the spread of the data. Expect to see the median, quartiles, and outliers for age and resting blood pressure.



- Scatter Matrix:

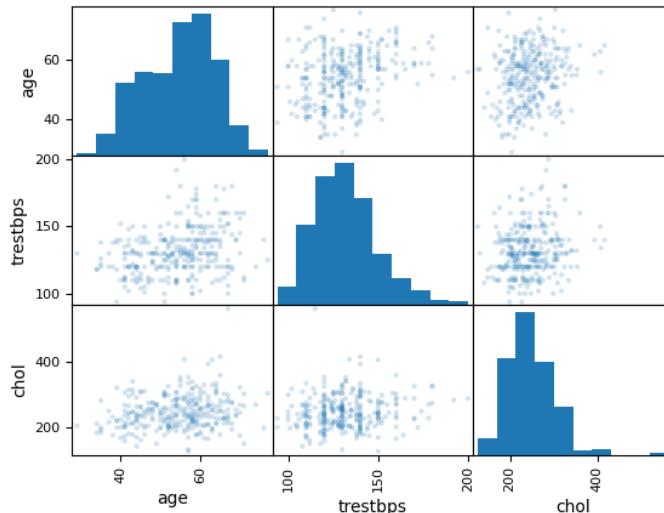
```

1 # Scatter matrix of multiple variables
2 pd.plotting.scatter_matrix(heart_disease_data[['age', 'trestbps', 'chol']],
3                           alpha=0.2, figsize=(10, 10))
4 plt.suptitle('Scatter Matrix of Age, Resting BP, and Cholesterol')
5 plt.show()

```

Description: Scatter matrices allow for the visualization of potential relationships between different numerical variables. Expect to see scatter plots for each pair of variables and histograms on the diagonal.

Scatter Matrix of Age, Resting BP, and Cholesterol



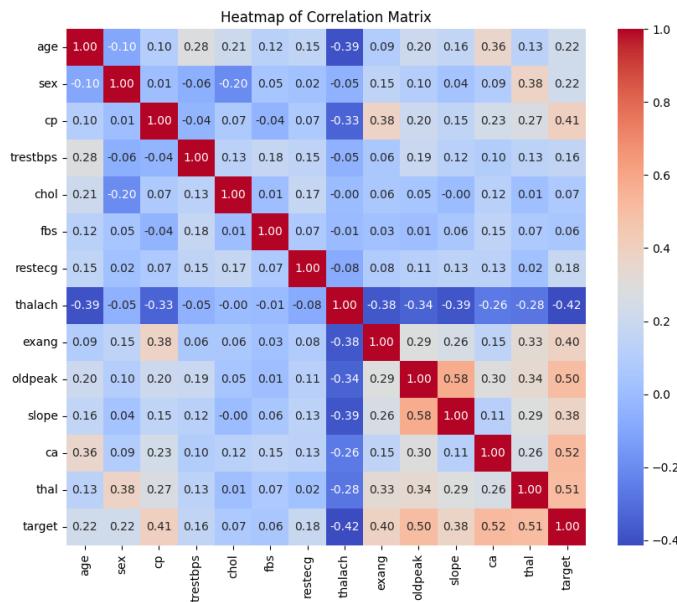
- Heatmap of Correlation Matrix:

```

1 # Heatmap showing correlations between variables
2 sns.heatmap(heart_disease_data.corr(), annot=True, cmap='coolwarm')
3 plt.title('Heatmap of Correlation Matrix')
4 plt.show()
5

```

Description: Heatmaps help in visualizing the strength of relationships between variables. Look for color-coded correlations, with warmer colors indicating stronger positive relationships and cooler colors indicating negative relationships.



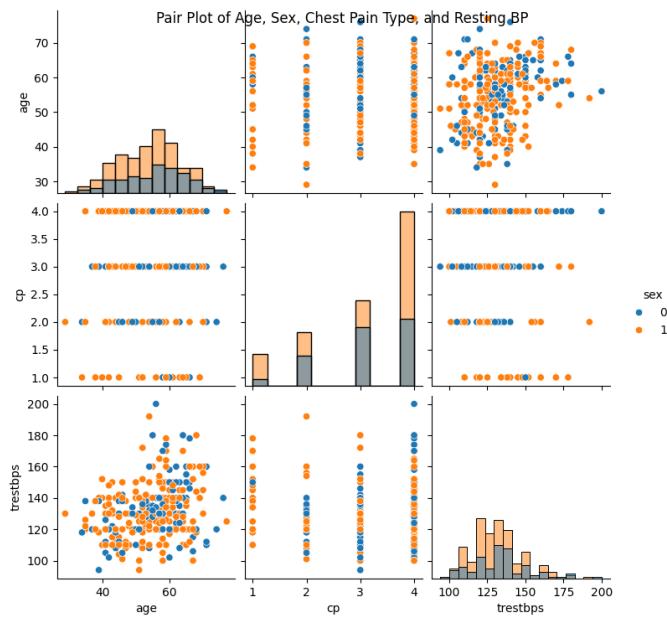
- **Pair Plot:**

```

1   # Pairwise plots of DataFrame's variables
2   sns.pairplot(heart_disease_data[['age', 'sex', 'cp', 'trestbps']], hue='sex')
3   plt.suptitle('Pair Plot of Age, Sex, Chest Pain Type, and Resting BP')
4   plt.show()
5

```

Description: Pair plots are useful for exploring correlations between multiple variables broken down by categories (colored by 'sex'). Expect to see a grid of plots showing relationships between pairs of variables.



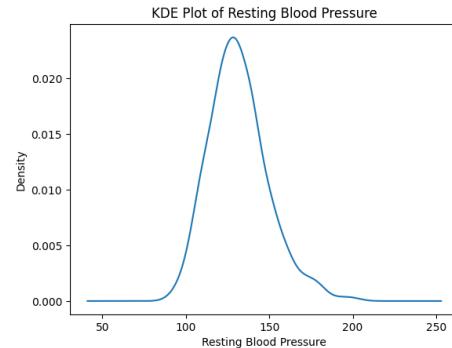
- Kernel Density Estimate Plot:

```

1 # KDE plot for 'trestbps'
2 heart_disease_data['trestbps'].plot(kind='kde')
3 plt.xlabel('Resting Blood Pressure')
4 plt.show()
5
6

```

Description: KDE plots are used to estimate the probability density function of a continuous random variable. Expect to see a smoothed version of the histogram, providing insights into the distribution of resting blood pressures.



5 Time Series Operations

Sample Data

Below is a sample of the weather dataset that will be used for the following examples. This data includes hourly recordings of various meteorological parameters.

Year	Month	Day	Hour	Temperature (°C)	Humidity (%)	MSL Pressure (hPa)	Wind Speed (km/h)
2018	11	6	0	21.65	38	1011.7	17.15
2018	11	6	1	20.94	40	1011.4	16.42
2018	11	6	2	20.29	41	1011.2	16.08
2018	11	6	3	19.67	43	1011	14.01
2018	11	6	4	18.76	45	1011.8	10.54
2018	11	6	5	18.25	47	1012.5	9.26
2018	11	6	6	17.68	48	1012.3	8.09
2018	11	6	7	19.63	42	1013.1	6.79

Handling Time Series

- Converting Columns to Datetime:

```

1 # Create a datetime index from the 'Year', 'Month', 'Day', 'Hour' columns
2 daily_weather['Datetime'] = pd.to_datetime(daily_weather[['Year', 'Month',
3   'Day', 'Hour']])
4 daily_weather.set_index('Datetime', inplace=True)

```

Description: Converting these columns into a single datetime column enables effective indexing and manipulation of time series data, essential for time-based analysis.

- Resampling Time-Series Data:

```

1 # Resample the data by month and calculate mean temperature
2 monthly_avg_temperature = daily_weather['Temperature'].resample('M').mean()
3

```

Description: Resampling helps summarize daily temperature data into monthly averages, simplifying the data and highlighting longer-term trends.

Window Functions

- Moving Average:

```

1 # Compute a moving average of the temperature over a 5-day window
2 daily_weather['5_day_temp_avg'] = daily_weather['Temperature'].rolling(
3   window=5).mean()

```

Description: A moving average smooths out short-term fluctuations and helps to identify longer-term trends in temperature data.

- Cumulative Sum:

```
1      # Calculate a cumulative sum of wind speed
2  daily_weather['cumulative_wind_speed'] = daily_weather['Wind Speed'].  
3  expanding(min_periods=1).sum()
```

Description: The cumulative sum of wind speed can reveal the total amount of wind experienced over a period, useful for climatological studies and analysis.

6 Matplotlib Visualization

Matplotlib is one of the most popular Python libraries for producing static, animated, and interactive visualizations in Python. This section will cover the basics of creating visualizations using Matplotlib, including line plots, bar charts, histograms, scatter plots, and pie charts.

Line Plots

Line plots are one of the simplest types of plots, useful for displaying trends over time.

```
1 import matplotlib.pyplot as plt
2
3 # Sample data
4 x = range(10)
5 y = [xi**2 for xi in x]
6
7 # Creating a line plot
8 plt.figure()
9 plt.plot(x, y)
10 plt.title('Simple Line Plot')
11 plt.xlabel('X Axis Label')
12 plt.ylabel('Y Axis Label')
13 plt.show()
```

Bar Charts

Bar charts are great for comparing numerical values across different categories.

```
1 # Sample data
2 categories = ['Category 1', 'Category 2', 'Category 3']
3 values = [10, 20, 15]
4
5 # Creating a bar chart
6 plt.figure()
7 plt.bar(categories, values)
8 plt.title('Simple Bar Chart')
9 plt.xlabel('Categories')
10 plt.ylabel('Values')
11 plt.show()
```

Histograms

Histograms are used to summarize the distribution of a dataset.

```
1 import numpy as np
2
3 # Generating random data
4 data = np.random.randn(1000)
5
6 # Creating a histogram
7 plt.figure()
8 plt.hist(data, bins=30)
9 plt.title('Simple Histogram')
10 plt.xlabel('Data Points')
11 plt.ylabel('Frequency')
12 plt.show()
```

Scatter Plots

Scatter plots are used to determine the relationship between two variables.

```
1 # Generating sample data
2 x = np.random.rand(50)
3 y = np.random.rand(50)
4
5 # Creating a scatter plot
6 plt.figure()
7 plt.scatter(x, y)
8 plt.title('Simple Scatter Plot')
9 plt.xlabel('X Axis')
10 plt.ylabel('Y Axis')
11 plt.show()
```

Pie Charts

Pie charts show the proportions of categories as parts of a whole.

```
1 # Sample data
2 labels = ['Part 1', 'Part 2', 'Part 3', 'Part 4']
3 sizes = [15, 30, 45, 10]
4
5 # Creating a pie chart
6 plt.figure()
7 plt.pie(sizes, labels=labels, autopct='%1.1f%%')
8 plt.axis('equal') # Equal aspect ratio ensures that pie is drawn as a circle.
9 plt.title('Simple Pie Chart')
10 plt.show()
```

7 Seaborn Visualization

Seaborn is a Python data visualization library based on Matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics. This section covers the basics of creating visualizations using Seaborn, including histograms, bar plots, count plots, box plots, and violin plots.

Histograms (Distplots)

Seaborn's distplot lets you show a histogram with a line on it. This can be used to display the distribution of the data.

```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Generating random data
5 data = np.random.randn(1000)
6
7 # Creating a distplot
8 plt.figure()
9 sns.distplot(data, bins=30, kde=True)
10 plt.title('Seaborn Histogram with Density Plot')
11 plt.show()
```

Bar Plots

Bar plots help you visualize the distributions of categorical data.

```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Sample data
5 categories = ['Category 1', 'Category 2', 'Category 3']
6 values = [10, 20, 15]
7
8 # Creating a bar plot
9 plt.figure()
10 sns.barplot(x=categories, y=values)
11 plt.title('Seaborn Bar Plot')
12 plt.xlabel('Categories')
13 plt.ylabel('Values')
14 plt.show()
```

Count Plots

Count plots can be thought of as histograms across a categorical, instead of quantitative, variable.

```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Sample categorical data
5 data = ['Category 1', 'Category 2', 'Category 1', 'Category 3', 'Category 2',
       'Category 3', 'Category 1']
6
7 # Creating a count plot
8 plt.figure()
9 sns.countplot(data)
10 plt.title('Seaborn Count Plot')
11 plt.xlabel('Categories')
12 plt.ylabel('Frequency')
13 plt.show()
```

Box Plots

Box plots show distributions with quartiles and whiskers. They can be used to compare distributions across the categories.

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 # Sample data
6 data = np.random.rand(10, 4)
7 df = pd.DataFrame(data, columns=['A', 'B', 'C', 'D'])
8
9 # Creating a box plot
10 plt.figure()
11 sns.boxplot(data=df)
12 plt.title('Seaborn Box Plot')
13 plt.show()
```

Violin Plots

Violin plots are used to visualize the distribution of the data and its probability density.

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 # Sample data
6 data = np.random.rand(10, 4)
7 df = pd.DataFrame(data, columns=['A', 'B', 'C', 'D'])
8
9 # Creating a violin plot
10 plt.figure()
11 sns.violinplot(data=df)
12 plt.title('Seaborn Violin Plot')
13 plt.show()
```