

# Python packages

Shucheng Liao\*

May 2024

## Contents

<b>1</b>	<b>Numpy</b>	<b>2</b>
<b>2</b>	<b>Scipy</b>	<b>4</b>
<b>3</b>	<b>Pandas Dataframe</b>	<b>7</b>
<b>4</b>	<b>Requests</b>	<b>10</b>
<b>5</b>	<b>Beautifulsoup4</b>	<b>11</b>
<b>6</b>	<b>Matplotlib</b>	<b>14</b>
<b>7</b>	<b>Seaborn</b>	<b>17</b>
<b>8</b>	<b>Scikit-learn</b>	<b>20</b>
<b>9</b>	<b>Thank You Note</b>	<b>23</b>

---

\*NYU Stern, MS in Quantitative Economics

# 1 Numpy

```
1 import numpy as np
2
3 ## Create a NumPy array
4 arr = np.array([1, 2, 3, 4, 5])
5 print(arr)
6 # Expected output: [1 2 3 4 5]
7
8 ## Create an array of zeros/ones
9 zeros = np.zeros((2, 3)) # np.ones((2,3))
10 print(zeros)
11 # Expected output:
12 # [[0. 0. 0.]
13 #  [0. 0. 0.]]
14
15 ## Create an array with a range of values
16 range_arr = np.arange(10)
17 print(range_arr)
18 # Expected output: [0 1 2 3 4 5 6 7 8 9]
19
20 ## Create an array with evenly spaced values
21 linspace_arr = np.linspace(0, 1, 5)
22 print(linspace_arr)
23 # Expected output: [0.    0.25 0.5   0.75 1.   ]
24
25 ## Reshape an array
26 reshaped = np.reshape(range_arr, (2, 5))
27 print(reshaped)
28 # Expected output:
29 # [[0 1 2 3 4]
30 #  [5 6 7 8 9]]
31
32 ## Basic arithmetic operations
33 a = np.array([1, 2, 3])
34 b = np.array([4, 5, 6])
35
36 print(a + b)
37 # Expected output: [5 7 9]
38
39 print(a - b)
40 # Expected output: [-3 -3 -3]
41
42 print(a * b)
43 # Expected output: [4 10 18]
44
45 print(a / b)
46 # Expected output: [0.25 0.4   0.5 ]
47
48 ## Element-wise square root
49 sqrt_arr = np.sqrt(a)
50 print(sqrt_arr)
51 # Expected output: [1.          1.41421356 1.73205081]
52
53 ## Dot product of two arrays
54 dot_product = np.dot(a, b)
55 print(dot_product)
56 # Expected output: 32
57
58 ## Sum of all elements
59 sum_all = np.sum(a)
60 print(sum_all)
61 # Expected output: 6
62
63 ## Sum along an axis
64 matrix = np.array([[1, 2], [3, 4]])
65 sum_axis0 = np.sum(matrix, axis=0)
```

```
66 | print(sum_axis0)
67 | # Expected output: [4 6]
68 |
69 | sum_axis1 = np.sum(matrix, axis=1)
70 | print(sum_axis1)
71 | # Expected output: [3 7]
72 |
73 | ## Transpose of an array
74 | transpose = np.transpose(matrix)
75 | print(transpose)
76 | # Expected output:
77 | # [[1 3]
78 | #   [2 4]]
79 |
80 | ## Boolean indexing
81 | bool_index = a > 2
82 | print(bool_index)
83 | # Expected output: [False False  True]
84 |
85 | filtered = a[a > 2]
86 | print(filtered)
87 | # Expected output: [3]
```

[Back to Table of Contents](#)

## 2 Scipy

### 1. Optimization

```
1 from scipy.optimize import minimize
2
3 # Define a quadratic function
4 def func(x):
5     return x**2 + 2*x + 1
6
7 # Minimize the function
8 result = minimize(func, 0) # initial guess 0
9 print(result.x)
10 # Expected output: [-1.]
```

Finds the minimum of the quadratic function, which is at  $x = -1$ .

### 2. Integration

```
1 from scipy.integrate import quad
2
3 # Define a simple function
4 def integrand(x):
5     return x**2
6
7 # Integrate the function from 0 to 1
8 result, error = quad(integrand, 0, 1)
9 print(result)
10 print(error)
11 # Expected output:
12 # result: 0.3333333333333333
13 # error: 3.700743415417189e-15
```

Computes the integral of  $x^2$  from 0 to 1, which is  $1/3$ . The error term gives an estimate of the numerical integration error, which is very small in this case.

### 3. Interpolation

```
1 from scipy.interpolate import interp1d
2 import numpy as np
3
4 # Define data points
5 x = np.array([0, 1, 2, 3])
6 y = np.array([0, 1, 0, 1])
7
8 # 1. Linear Interpolation
9 f_linear = interp1d(x, y, kind='linear')
10 print(f_linear(1.5)) # estimate the value of y at x = 1.5
11 # Expected output: 0.5
```

Linear Interpolation: Connects data points with straight lines, simple and fast. Commonly used for quick approximations.

```
1 # 2. Nearest-Neighbor Interpolation
2 f_nearest = interp1d(x, y, kind='nearest')
3 print(f_nearest(1.5))
4 # Expected output: 0.0
```

Nearest-Neighbor Interpolation: Uses the nearest data point's value. Useful for categorical data or when a step function is desired.

```
1 # 3. Zero-Order Hold (Previous)
2 f_zero = interp1d(x, y, kind='zero')
```

```

3 print(f_zero(1.5))
4 # Expected output: 1.0

```

Zero-Order Hold: Uses the previous data point's value (step function). Suitable for systems where the value remains constant until the next data point.

```

1 # 4. Spline Interpolation of Order 1
2 f_slinear = interp1d(x, y, kind='slinear')
3 print(f_slinear(1.5))
4 # Expected output: 0.5

```

Spline Interpolation of Order 1: Similar to linear interpolation but uses splines for a smoother transition. It can handle more complex data structures better than simple linear interpolation in some cases.

```

1 # 5. Quadratic/Cubic Spline Interpolation
2 f_quadratic = interp1d(x, y, kind='quadratic') # 'cubic'
3 print(f_quadratic(1.5))
4 # Expected output: Depends on the quadratic fit

```

Quadratic/Cubic Spline Interpolation: Fits a quadratic/cubic polynomial between data points.

## 4. Linear Algebra

```

1 from scipy.linalg import eig
2 import numpy as np
3
4 # Define a matrix
5 A = np.array([[1, 2], [3, 4]])
6
7 # Compute eigenvalues and eigenvectors
8 values, vectors = eig(A)
9 print(values)
10 # Expected output: Array of eigenvalues [5.37228132 -0.37228132]
11 print(vectors)
12 # Expected output: Matrix of eigenvectors
13 # [[-0.41597356 -0.82456484]
14 #  [-0.90937671  0.56576746]]

```

Computes the eigenvalues and eigenvectors of matrix A.

## 5. Statistical Functions

```

1 from scipy.stats import norm
2
3 # Define a normal distribution
4 mean, std_dev = 0, 1
5 dist = norm(mean, std_dev)
6
7 # Compute the cumulative distribution function at 0
8 print(dist.cdf(0)) # (P(X <= 0)) # (dist.logcdf(0))
9 # Expected output: 0.5
10
11 # Compute the probability density function at 0
12 print(dist.pdf(0)) # (dist.logpdf(0))
13 # Expected output: 0.3989422804014327

```

Computes the CDF/PDF of the standard normal distribution at  $x = 0$ .

## 7. Spatial Data

```
1 from scipy.spatial import distance
2
3 # Define two points
4 point1 = [0, 0]
5 point2 = [3, 4]
6
7 # Compute Euclidean distance
8 dist = distance.euclidean(point1, point2)
9 print(dist)
10 # Expected output: 5.0
```

Computes the Euclidean distance between two points.

[Back to Table of Contents](#)

## 3 Pandas Dataframe

### 1. Reading Data

```
1 import pandas as pd
2
3 # Reading data from a CSV file
4 df = pd.read_csv('economic_data.csv')
5 # or pd.read_excel or pd.read_stata
6 print(df.head())      # First 5 rows of the CSV file
```

### 2. Inspecting Data

```
1 # Inspecting the DataFrame
2 print(df.info())
3 # Expected output: Summary of the DataFrame, including column names, non-null
  counts, and data types
4
5 print(df.describe())
6 # Expected output: Descriptive statistics for numerical columns: count, mean, std,
  min, max, percentiles
7
8 # Descriptive statistics for categorical data
9 print(df.describe(include=['object']))
10 # Expected output: count, unique, top, and frequency for categorical columns
```

### 3. Selecting Data

```
1 # Selecting specific column
2 gdp_series = df['GDP']
3 print(gdp_series.head())
4 # Expected output: First 5 values in the 'GDP' column
5
6 # Selecting specific subset
7 subset = df[['GDP', 'Inflation']]
8 print(subset.head())
9 # Expected output: First 5 rows of the DataFrame with 'GDP' and 'Inflation' columns
```

### 4. Filtering Data

```
1 # Filtering rows based on a condition
2 filtered_df = df[df['GDP'] > 10000]
3 print(filtered_df.head())
4 # Expected output: Rows where GDP is greater than 10,000
```

### 5. Adding New Columns

```
1 # Adding a new column
2 df['GDP per Capita'] = df['GDP'] / df['Population']
3 print(df.head())
4 # Expected output: DataFrame with a new 'GDP per Capita' column
```

## 6. Grouping Data

```
1 # Grouping data
2 grouped_df = df.groupby('Country')['GDP'].sum()
3 print(grouped_df)
4 # Expected output: Sum of GDP for each country
5
6 # Grouping data and applying an aggregation function
7 grouped_df = df.groupby('Country')['GDP'].agg(['mean', 'sum', 'min', 'max'])
8 # Expected output: DataFrame with mean, sum, min, and max GDP for each country
```

## 7. Merging DataFrames

```
1 # Merging two DataFrames
2 df1 = pd.DataFrame({'Country': ['A', 'B'], 'GDP': [1000, 2000]})
3 df2 = pd.DataFrame({'Country': ['A', 'B'], 'Population': [10, 20]})
4 merged_df = pd.merge(df1, df2, on='Country')
5
6 # Extension: how='inner'/'outer'/'left'/'right' - includes rows with matched keys
7 # only/ all rows with unmatched keys marked as 'NaN'/ all rows from left dt and
8 # matched rows from right dt/ all rows from right dt and matched rows from left dt
9
10 print(merged_df)
11 # Expected output: Merged DataFrame with GDP and Population for each country
```

## 8. Handling Missing Data

```
1 # Fill missing values with a specific value or method
2 merged_df.fillna(0, inplace=True)
3
4 # Drop rows with any missing values
5 merged_df.dropna(inplace=True)
6
7 # Handling missing data - filled with the mean GDP
8 df['GDP'].fillna(df['GDP'].mean(), inplace=True)
9 print(df.head())
```

## 9. Time Series Analysis

```
1 # form of date to be convert (example): ['01-01-2024', '15-02-2024']
2
3 # Converting a column to datetime
4 df['Date'] = pd.to_datetime(df['Date'])
5
6 # Extension: format='%d-%m-%Y' - switching the orders
7
8 # Extension: errors = 'coerce'/'raise'/'ignore' - convert the invalid to NaT/raise an
9 # error/return original if fail to convert
10
11 # Setting it as the index
12 df.set_index('Date', inplace=True)
13 print(df.head())
14 # Expected output: DataFrame with 'Date' column as the index
```

## 10. Resampling Time Series Data



```
1 # Resampling time series data to a different frequency
2 monthly_df = df.resample('M').mean()
3 print(monthly_df.head())
4 # Expected output: DataFrame resampled to monthly frequency with mean values
```

## 11. Exporting DataFrame to Different File Formats

```
1 # Export to CSV
2 df.to_csv('output.csv', index=False)
3 # Saves the DataFrame to a CSV file named 'output.csv' without the index.
4
5 # Extension: to_excel/to_stata
```

[Back to Table of Contents](#)

## 4 Requests

```
1 # Sending a GET request and checking for errors
2 response = requests.get('https://api.example.com/data')
3 try:
4     response.raise_for_status()
5     print('Request was successful')
6     # Expected output: 'Request was successful' if the status code indicates success
7 except requests.HTTPError as e:
8     print(f'HTTP error occurred: {e}')
9     # Expected output: 'HTTP error occurred: ...' with details of the error
```

[Back to Table of Contents](#)

## 5 BeautifulSoup4

### Quick note: Extract text from HTML requests

```
1 import requests
2 from bs4 import BeautifulSoup
3
4 # Send a GET request
5 response = requests.get('https://www.example.com')
6
7 # Check if the request was successful
8 if response.status_code == 200:
9     # Parse the HTML content
10    soup = BeautifulSoup(response.content, 'html.parser')
11
12    # Extract text from the HTML content
13    text = soup.get_text()
14
15    print(text)
16    # Expected output: All text content extracted from the HTML page
17 else:
18    print(f'Failed to retrieve data: {response.status_code}')
19    # Expected output: Failure message with HTTP status code
```

### 1. Creating a BeautifulSoup Object

```
1 from bs4 import BeautifulSoup
2
3 # Sample HTML content
4 html_content = '''
5 <html>
6     <head><title>Economic Data</title></head>
7     <body>
8         <h1>Economic Indicators</h1>
9         <p id="gdp">GDP: $20 Trillion</p>
10        <p id="inflation">Inflation: 2%</p>
11    </body>
12 </html>
13 '''
14
15 # Creating a BeautifulSoup object
16 soup = BeautifulSoup(html_content, 'html.parser')
17 print(soup.prettify())
18 # Expected output: Formatted HTML content
```

### 2. Accessing Elements by Tag Name

```
1 # Accessing elements by tag name
2 title = soup.title
3 print(title)
4 # Expected output: <title>Economic Data</title>
5
6 print(title.text)
7 # Expected output: Economic Data
```

### 3. Finding Elements by ID

```
1 # Finding elements by ID
2 gdp = soup.find(id='gdp')
3 print(gdp)
```

```

4 # Expected output: <p id="gdp">GDP: $20 Trillion</p>
5
6 print(gdp.text)
7 # Expected output: GDP: $20 Trillion

```

## 4. Finding Elements by Class Name

```

1 # Sample HTML content with class names
2 html_content = '''
3 <html>
4     <body>
5         <div class="economic-indicator">GDP: $20 Trillion</div>
6         <div class="economic-indicator">Inflation: 2%</div>
7     </body>
8 </html>
9 '''
10
11 # Creating a BeautifulSoup object
12 soup = BeautifulSoup(html_content, 'html.parser')
13
14 # Finding elements by class name
15 indicators = soup.find_all(class_='economic-indicator')
16 for indicator in indicators:
17     print(indicator.text)
18 # Expected output:
19 # GDP: $20 Trillion
20 # Inflation: 2%

```

## 5. Extracting Links (Anchor Tags)

```

1 # Sample HTML content with links
2 html_content = '''
3 <html>
4     <body>
5         <a href="https://example.com/gdp">GDP Report</a>
6         <a href="https://example.com/inflation">Inflation Report</a>
7     </body>
8 </html>
9 '''
10
11 # Creating a BeautifulSoup object
12 soup = BeautifulSoup(html_content, 'html.parser')
13
14 # Extracting all links
15 links = soup.find_all('a')
16 for link in links:
17     print(link.get('href'))
18 # Expected output:
19 # https://example.com/gdp
20 # https://example.com/inflation

```

## 6. Navigating the Parse Tree

```

1 header = soup.body.h1
2 print(header)
3 # Expected output: <h1>Economic Indicators</h1>
4
5 paragraphs = soup.body.find_all('p')
6 for para in paragraphs:
7     print(para.text)
8 # Expected output:

```

```
9 # GDP: $20 Trillion
10 # Inflation: 2%
```

## 7. Modifying the HTML Content

```
1 soup.body.h1.string = "Updated Economic Indicators"
2 print(soup.body.h1)
3 # Expected output: <h1>Updated Economic Indicators</h1>
4
5 new_tag = soup.new_tag("p", id="unemployment")
6 new_tag.string = "Unemployment: 5%"
7 soup.body.append(new_tag)
8 print(soup.body)
9 # Expected output: <body>...<p id="unemployment">Unemployment: 5%</p></body>
```

## 8. Removing Elements

```
1 for tag in soup.find_all('p'):
2     tag.decompose()
3
4 print(soup.body)
5 # Expected output: <body>...</body> without <p> tags
```

[Back to Table of Contents](#)

## 6 Matplotlib

see different shapes of markers

see different named colors

see different line styles

### 1. Basic Line Plot

```
1 import matplotlib.pyplot as plt
2
3 # Sample economic data
4 years = [2016, 2017, 2018, 2019, 2020]
5 gdp = [1.5, 2.3, 2.9, 2.3, -3.5]
6
7 # Creating a basic line plot
8 plt.plot(years, gdp, marker='o')
9 plt.title('GDP Growth Over Years')
10 plt.xlabel('Year')
11 plt.ylabel('GDP Growth (%)')
12 plt.grid(True)
13 plt.show()
14
15 # Expected output: a basic line plot of GDP growth over the years with markers at
    data points
```

### 2. Bar Chart

```
1 import matplotlib.pyplot as plt
2
3 # Sample economic data
4 countries = ['USA', 'China', 'Japan', 'Germany', 'India']
5 gdp = [21.43, 14.34, 5.08, 3.84, 2.87]
6
7 # Creating a bar chart
8 plt.bar(countries, gdp, color='skyblue')
9 plt.title('GDP of Top 5 Economies in 2020')
10 plt.xlabel('Country')
11 plt.ylabel('GDP (Trillions USD)')
12 plt.show()
13
14 # Expected output: Bar chart showing GDP of the top 5 economies in 2020
```

### 3. Pie Chart

```
1 import matplotlib.pyplot as plt
2
3 # Sample economic data
4 sectors = ['Agriculture', 'Industry', 'Services']
5 gdp_share = [4, 33, 63]
6
7 # Creating a pie chart
8 plt.pie(gdp_share, labels=sectors, autopct='%1.1f%%', startangle=90)
9 plt.title('GDP Share by Sector')
10 plt.show()
11
12 # Expected output: Pie chart showing GDP share by sector
13 # '%1.0f%%'/'%1.1f%%'/'%1.2f%%': 0/1/2 decimal place(s)
14 # startangle=90: This will start the first pie slice from the top
```

## 4. Scatter Plot

```
1 import matplotlib.pyplot as plt
2
3 # Sample economic data
4 education_spending = [4.1, 3.5, 3.8, 3.2, 3.0]
5 gdp_growth = [2.9, 1.6, 2.2, 2.0, 1.8]
6
7 # Creating a scatter plot
8 plt.scatter(education_spending, gdp_growth, color='red')
9 plt.title('Education Spending vs GDP Growth')
10 plt.xlabel('Education Spending (% of GDP)')
11 plt.ylabel('GDP Growth (%)')
12 plt.grid(True)
13 plt.show()
14
15 # Expected output: Scatter plot showing the relationship between education spending
    and GDP growth
```

## 5. Histogram

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 # Sample economic data
5 inflation_rates = np.random.normal(3, 1, 1000) # Generate 1000 random data points
6
7 # Creating a histogram
8 plt.hist(inflation_rates, bins=20, color='green', edgecolor='black')
9 plt.title('Inflation Rate Distribution')
10 plt.xlabel('Inflation Rate (%)')
11 plt.ylabel('Frequency')
12 plt.show()
13 # Expected output: Histogram showing the distribution of inflation rates
```

## 6. Time Series Plot

```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3
4 # Sample time series data
5 date_range = pd.date_range(start='1/1/2020', periods=12, freq='M')
6 unemployment_rate = [5.1, 5.2, 5.0, 4.9, 5.3, 5.5, 5.6, 5.4, 5.3, 5.2, 5.1, 5.0]
7
8 # Creating a time series plot
9 plt.plot(date_range, unemployment_rate, marker='o', linestyle='--')
10 plt.title('Monthly Unemployment Rate in 2020')
11 plt.xlabel('Month')
12 plt.ylabel('Unemployment Rate (%)')
13 plt.grid(True)
14 plt.show()
15 # Expected output: Time series plot showing the monthly unemployment rate in 2020
```

## 7. Adding Annotations

```
1 import matplotlib.pyplot as plt
2
3 # Sample economic data
4 years = [2016, 2017, 2018, 2019, 2020]
5 gdp = [1.5, 2.3, 2.9, 2.3, -3.5]
6
```

```

7 # Creating a line plot with annotations
8 plt.plot(years, gdp, marker='o')
9 plt.title('GDP Growth Over Years')
10 plt.xlabel('Year')
11 plt.ylabel('GDP Growth (%)')
12 plt.grid(True)
13
14 # Adding annotation
15 plt.annotate('COVID-19 Impact', xy=(2020, -3.5), xytext=(2018, 0),
16             arrowprops=dict(facecolor='black', shrink=0.05))
17 plt.show()
18 # Expected output: Line plot with an annotation highlighting the impact of COVID-19
    in 2020

```

## 8. Adding More Stuff

```

1 import matplotlib.pyplot as plt
2
3 # Sample economic data
4 years = [2016, 2017, 2018, 2019, 2020]
5 gdp = [1.5, 2.3, 2.9, 2.3, -3.5]
6
7 # Creating a basic line plot
8 plt.plot(years, gdp, marker='o', label='GDP Growth')
9
10 # Adding a horizontal line at y=0 with a dashed line style
11 plt.axhline(y=0, color='gray', linestyle='--')
12
13 # Adding a vertical line at x=2018 with a dotted line style
14 plt.axvline(x=2018, color='blue', linestyle=':')
15
16 # Adding a marker at a specific coordinate (2019, 2.3)
17 plt.plot(2019, 2.3, marker='x', markersize=10, color='red')
18
19 # Adding a 45-degree line (y=x)
20 x_values = range(2016, 2021)
21 plt.plot(x_values, x_values, color='green', linestyle='-.', label='45-degree line')
22
23 # Adding a legend
24 plt.legend()
25
26 plt.title('GDP Growth Over Years with Enhancements')
27 plt.xlabel('Year')
28 plt.ylabel('GDP Growth (%)')
29 plt.grid(True)
30 plt.show()

```

[Back to Table of Contents](#)



## 7 Seaborn

see different color palettes

### 1. Importing Seaborn and Matplotlib

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
```

### 2. Setting Aesthetic Style

```
1 sns.set(style="whitegrid")
2 # Sets the aesthetic style of the plots to whitegrid, which is useful for economic
  data visualizations.
```

### 3. Line Plot

```
1 # Sample economic data
2 years = [2016, 2017, 2018, 2019, 2020]
3 gdp = [1.5, 2.3, 2.9, 2.3, -3.5]
4
5 # Creating a line plot
6 sns.lineplot(x=years, y=gdp)
7 plt.title('GDP Growth Over Years')
8 plt.xlabel('Year')
9 plt.ylabel('GDP Growth (%)')
10 plt.show()
11 # Expected output: Line plot showing GDP growth over the years.
```

### 4. Bar Plot

```
1 # Sample economic data
2 countries = ['USA', 'China', 'Japan', 'Germany', 'India']
3 gdp = [21.43, 14.34, 5.08, 3.84, 2.87]
4
5 # Creating a bar plot
6 sns.barplot(x=countries, y=gdp, palette="Blues_d")
7 plt.title('GDP of Top 5 Economies in 2020')
8 plt.xlabel('Country')
9 plt.ylabel('GDP (Trillions USD)')
10 plt.show()
11 # Expected output: Bar plot showing GDP of the top 5 economies in 2020.
```

### 5. Histogram

```
1 import numpy as np
2
3 # Sample economic data
4 inflation_rates = np.random.normal(3, 1, 1000) # Generate 1000 random data points
5
6 # Creating a histogram
7 sns.histplot(inflation_rates, bins=20, kde=True, color="green")
8 plt.title('Inflation Rate Distribution')
9 plt.xlabel('Inflation Rate (%)')
10 plt.ylabel('Frequency')
11 plt.show()
```

```

12 # Expected output: Histogram showing the distribution of inflation rates with KDE:
    this parameter adds a KDE plot over the histogram, providing a smooth estimate
    of the distribution of the data, especially when dealing with continuous data

```

## 6. Scatter Plot

```

1 # Sample economic data
2 education_spending = [4.1, 3.5, 3.8, 3.2, 3.0]
3 gdp_growth = [2.9, 1.6, 2.2, 2.0, 1.8]
4
5 # Creating a scatter plot
6 sns.scatterplot(x=education_spending, y=gdp_growth, color="red")
7 plt.title('Education Spending vs GDP Growth')
8 plt.xlabel('Education Spending (% of GDP)')
9 plt.ylabel('GDP Growth (%)')
10 plt.show()
11 # Expected output: Scatter plot showing the relationship between education spending
    and GDP growth.

```

## 7. Pair Plot

```

1 # Sample economic data
2 data = {
3     'GDP': [1.5, 2.3, 2.9, 2.3, -3.5],
4     'Inflation': [2.1, 1.8, 2.5, 1.9, 0.5],
5     'Unemployment': [5.1, 4.9, 4.7, 5.2, 6.1]
6 }
7 df = pd.DataFrame(data, index=[2016, 2017, 2018, 2019, 2020])
8
9 # Creating a pair plot
10 sns.pairplot(df)
11 plt.show()
12 # Expected output: Pair plot showing relationships between GDP, inflation, and
    unemployment.

```

## 8. Heatmap

```

1 # Sample economic data
2 data = np.array([[1.5, 2.3, 2.9], [2.1, 1.8, 2.5], [5.1, 4.9, 4.7]])
3 years = ['2018', '2019', '2020']
4 indicators = ['GDP Growth', 'Inflation', 'Unemployment']
5
6 # Creating a heatmap
7 sns.heatmap(data, annot=True, xticklabels=indicators, yticklabels=years,
8             cmap="YlGnBu")
9 plt.title('Economic Indicators Heatmap')
10 plt.show()
11 # Expected output: Heatmap showing economic indicators over different years.
12 # cmap="YlGnBu"- Sequential Colormaps representing data that has a natural order
    (e.g., low to high values).

```

## 9. Box Plot

```

1 # Sample economic data
2 data = {
3     'Country': ['USA', 'China', 'Japan', 'Germany', 'India'] * 5,
4     'Year': [2016, 2017, 2018, 2019, 2020] * 5,
5     'GDP Growth': [1.5, 2.3, 2.9, 2.3, -3.5, 6.7, 6.9, 6.6, 6.1, 2.3, 1.0, 0.3, 0.8,
6                   0.5, 1.8, 2.2, 2.5, 2.0, 2.1, 1.1, 7.1, 6.9, 6.5, 6.0, 5.5]

```

```
6 }
7 df = pd.DataFrame(data)
8
9 # Creating a box plot
10 sns.boxplot(x='Year', y='GDP Growth', data=df, palette="Set3")
11 plt.title('GDP Growth Distribution by Year')
12 plt.xlabel('Year')
13 plt.ylabel('GDP Growth (%)')
14 plt.show()
15 # Expected output: Box plot showing GDP growth distribution by year.
```

[Back to Table of Contents](#)

## 8 Scikit-learn

### 1. Importing Scikit-Learn Modules

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import mean_squared_error, r2_score
```

### 2. Loading and Splitting Data

```
1 from sklearn.datasets import load_boston
2
3 # Load dataset
4 data = load_boston()
5 X = data.data # feature matrix (independent variables, array/pd-dtf)
6 y = data.target # target vector (dependent variable, array/pd-dtf)
7
8 # Split dataset into training and test sets
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
10 random_state=42)
11
12 # test_size: proportion of the dataset to include in the test split, and the
13 # remaining will be used for training (use interger for absolute number of test
14 # samples)
15
16 # random_state: Using a fixed number allows you to get the same split every time you
17 # run the code: ensuring reproducibility of the split.
18
19 # Extension: shuffle=True/False - shuffle the data before splitting
```

### 3. Standardizing Data

```
1 # Standardize features
2 # Create an instance of the StandardScaler
3 # The StandardScaler standardizes features by removing the mean and scaling to unit
4 # variance.
5
6 scaler = StandardScaler()
7
8 # Fit the StandardScaler to the training data (X_train) and transform it
9 # This calculates the mean and standard deviation for scaling using the training data
10 # and then scales the training data based on these calculations.
11
12 X_train_scaled = scaler.fit_transform(X_train)
13
14 # Transform the test data (X_test) using the same mean and standard deviation
15 # calculated from the training data
16 # This ensures that the test data is scaled in the same way as the training data.
17
18 X_test_scaled = scaler.transform(X_test)
```

### 4. Training a Linear Regression Model and Predict

```
1 # Train a linear regression model
2 model = LinearRegression()
3 model.fit(X_train_scaled, y_train)
4
5 # Make predictions on the test set
6 y_pred = model.predict(X_test_scaled)
```

```

7 print(y_prob)
8 # Expected output: Array of predicted probabilities for each class, e.g., [[0.7,
    0.3], [0.2, 0.8], ...]

```

## 5. Evaluating the Model

```

1 # Evaluate the model
2 mse = mean_squared_error(y_test, y_pred)
3 r2 = r2_score(y_test, y_pred)
4 print(f'Mean Squared Error: {mse}')
5 print(f'R^2 Score: {r2}')

```

## 6. Cross-Validation

```

1 from sklearn.model_selection import cross_val_score
2
3 # Perform cross-validation
4 cv_scores = cross_val_score(model, X, y, cv=5)
5 print(f'Cross-Validation Scores: {cv_scores}')
6 print(f'Average CV Score: {cv_scores.mean()}')

```

## 7. Hyperparameter Tuning with Grid Search

```

1 from sklearn.model_selection import GridSearchCV
2 from sklearn.ensemble import RandomForestRegressor
3
4 # Define the model and parameter grid
5 rf = RandomForestRegressor(random_state=42)
6 param_grid = {
7     'n_estimators': [50, 100, 200],
8     'max_depth': [None, 10, 20, 30]
9 }
10
11 # Perform grid search
12 grid_search = GridSearchCV(rf, param_grid, cv=5, scoring='neg_mean_squared_error')
13 grid_search.fit(X_train, y_train)
14
15 # Best parameters and score
16 print(f'Best Parameters: {grid_search.best_params()}')
17 print(f'Best Score: {grid_search.best_score_}')

```

## 8. Feature Importance

```

1 # Feature importance from RandomForestRegressor
2 best_rf = grid_search.best_estimator_
3 importances = best_rf.feature_importances_
4 for feature, importance in zip(data.feature_names, importances):
5     print(f'{feature}: {importance}')

```

## 9. Pipeline for Workflow Automation

```

1 from sklearn.pipeline import Pipeline
2
3 # Create a pipeline with scaling and model training
4 pipeline = Pipeline([
5     ('scaler', StandardScaler()),
6     ('regressor', LinearRegression())

```

```
7 |])
8 |
9 | # Train the pipeline
10 | pipeline.fit(X_train, y_train)
11 |
12 | # Make predictions and evaluate
13 | y_pred_pipeline = pipeline.predict(X_test)
14 | mse_pipeline = mean_squared_error(y_test, y_pred_pipeline)
15 | r2_pipeline = r2_score(y_test, y_pred_pipeline)
16 | print(f'Pipeline Mean Squared Error: {mse_pipeline}')
17 | print(f'Pipeline R^2 Score: {r2_pipeline}')
```

[Back to Table of Contents](#)

## **9 Thank You Note**

Thank you, ChatGPT 4o, for helping me writing this guide.