

Python Data Analysis & Visualization

Schedule

Week 1: Foundations of Python

Class 1 (Day 1):

- a) Introduction to Python programming.
- b) Setting up the Python environment.
- c) Basic Python syntax and data types.

Class 2 (Day 2):

- a) Advanced Python concepts: functions, libraries, and modules.

Week 2: Data Analysis with NumPy and Pandas

Class 3 (Day 3):

- a) Introduction to NumPy
- b) Introduction to pandas for data handling.
- c) Importing and exporting data with pandas.
- d) Data cleaning and preprocessing.

Class 4 (Day 4):

- a) Data aggregation and transformation.
- b) Basic data analysis tasks using pandas.

Week 3: Data Visualization and Real-World Projects

Class 5 (Day 4):

- a) Introduction to data visualization libraries (Matplotlib).
- b) Creating basic plots and charts.

Class 6 (Day 5):

- a) Advanced data visualization techniques (Matplotlib and Seaborn).
- b) Real-world data analysis and visualization projects.

Step by Step Documentation on Python Data Analysis & Visualization

Python is a popular high-level programming language created by Guido van Rossum and released in 1991. It is widely used for research, data analysis, web development, software development, mathematics, and system scripting. If you are completely new to programming, Python is considered easy for beginners to use and learn. It is also easy for experienced programmers in any programming language to pick up. Installing Python is generally easy, and many Linux and UNIX distributions include a recent Python. Even some Windows computers now come with Python already installed.

Python offers several advantages that contribute to its popularity and widespread use in various domains.

Some of the key advantages of using Python include:

1. **Simplicity and Readability:** Python is known for its simplicity and readability, making it an excellent choice for both beginners and experienced programmers.
2. **Versatility and Support for Multiple Paradigms:** It supports multiple programming paradigms including procedural, object-oriented, and functional programming, and is used in various organizations due to its versatility.
3. **Rich Ecosystem of Third-Party Modules and Libraries:** Python has a rich ecosystem of third-party modules and extensive support libraries, making it suitable for scientific, data-related applications, and more.
4. **Open Source and Large Active Community:** Python is open source and has a large active community, which contributes to its continuous development and support.
5. **Interpreted Language and Easier Debugging:** Being an interpreted language, Python allows for easier debugging and code development.
6. **Extensive Libraries and User-Friendly Data Structures:** Python provides extensive support libraries and user-friendly data structures, which contribute to increased productivity and speed.
7. **Easy to Learn and Fun to Use:** Python's syntax reads like English, making it easy to learn and fun to use, which is advantageous for both beginners and experienced programmers.

These advantages make Python a popular choice for various applications, including web development, data science, artificial intelligence, and more.

Effectively teaching Python with a focus on data analysis involves providing learners with a thorough grasp of fundamental concepts and essential practices to establish a strong foundation. Here's an extensive guide I've prepared. On the first day, we will cover the following topics.

1. Installation and Setup

Start your instruction by guiding students through the process of installing Python and configuring an optimal development environment. While introducing them to well-known Integrated Development Environments (IDEs) like VSCode, PyCharm, or Jupyter Notebook, consider recommending the [Anaconda](#) distribution as a robust option. Then, use [VsCode](#) / Spyder / Jupyter notebook.

Why Anaconda?

Anaconda offers several advantages for data analysis-centric Python programming. It comes bundled with an extensive library of pre-installed packages and tools, which significantly simplifies the setup process for data-related tasks. One particularly invaluable inclusion is Jupyter Notebook, an interactive web application that facilitates code development, documentation, and visualization within a single environment. Its seamless integration with Anaconda makes it an indispensable tool for data analysis enthusiasts.

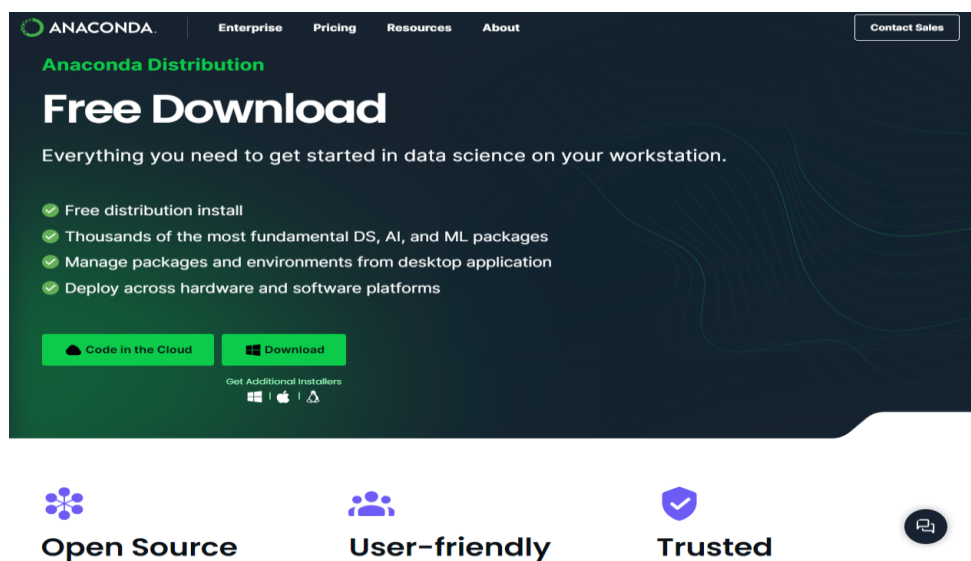


Figure 1. Anaconda Distribution webpage. It is available for Windows/ Mac/Linux OS.

i. **Installation Steps:**

Walk students through the installation process of Anaconda, ensuring they download the version compatible with their operating system. Emphasize the importance of selecting the appropriate Python version (**Python 3.11** is recommended for its compatibility with modern libraries).

Alternative: Google Colab: <https://colab.research.google.com/>

ii. **Environment Management:**

Teach students how to create and manage Python environments using Conda, a package and environment management tool included with Anaconda. Encourage the use of virtual environments to isolate projects and dependencies, preventing conflicts and ensuring project-specific package versions.

iii. **Package Management:**

Introduce **Conda** and **pip** as package managers to install, update, and remove Python packages and libraries.

Usage via conda / miniconda: **conda install anaconda::pandas**

Via pip : **pip install pandas**

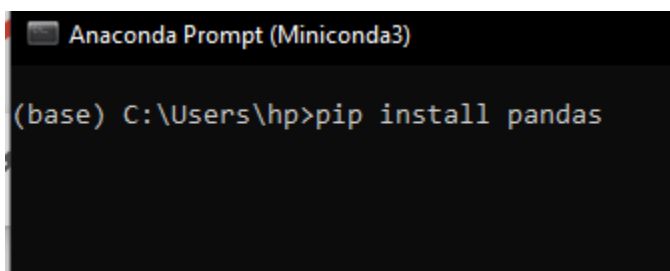


Figure 2. Package installation.

Explain the concept of requirements files for reproducible environments.

iv. **Jupyter Notebook:**

Dedicate time to exploring Jupyter Notebook, demonstrating its capabilities in code execution, visualization, and interactive data exploration. Show how to **create, run, export, and share** Jupyter

notebooks effectively. By equipping students with these skills and knowledge during the initial phase of your course, you lay a strong foundation for their journey into data analysis with Python.

Python

Python has 35 keywords, that it reserves for special syntactic uses. They are case sensitive. You cannot use keyword as regular identifiers.

Table 1. Python keywords

def	if	else	elif	for	in
class	import	and	return	True	False
or	while	try	None	del	yield
with	pass	break	continue	lambda	except

2. Basic Syntax:

In this crucial section of Python instruction, we delve into the very bedrock of the language – its syntax. The aim is to equip learners with a profound understanding of Python's syntax, encompassing core elements that form the building blocks of every Python program.

- i. **Indentation and Readability:** Python's elegance is further highlighted by its use of indentation for code blocks. We stress the significance of proper indentation and readability in Python, helping learners appreciate the importance of clean and organized code.
- ii. **Comments and Documentation:** Beyond the code itself, we introduce the practice of commenting code using # and writing docstrings to document functions and modules effectively. This instills in learners the habit of creating self-explanatory code.
- iii. **Code Structure:** We guide learners in understanding Python's code structure, where indentation and colons (:) are used to define blocks of code, such as loops and conditionals. This clarity of structure is a hallmark of Python programming.

- iv. **Variables:** We introduce learners to the concept of variables, explaining how they serve as placeholders for data and how to declare and assign values to them. We emphasize the dynamic nature of Python, where variables can change types during runtime.

```
# Variables and Data Types:
name = "John" # name is a variable of type string
age = 30 # age is a variable of type integer
height = 6.2 # height is a variable of type float
is_student = False # is_student is a variable of type boolean

print(name)
print(age)
print(height)
print(is_student)
```

- v. **Data Types:** Next, we explore Python's versatile data types. We cover the essentials, including integers (int), floating-point numbers (float), and strings (str). We emphasize the importance of data typing in Python and its impact on operations.

```
# show the type of each variable
print(type(name))
print(type(age))
print(type(height))
print(type(is_student))
```

- vi. **Basic Operations:** Python's power lies in its simplicity, and learners will become adept at using this simplicity to their advantage. We teach fundamental operations, such as arithmetic (addition, subtraction, multiplication, division), string manipulation (concatenation, slicing), and comparison (equality, inequality).

Table 2. Operators.

+	Addition		Bitwise OR
-	Subtraction	^	Bitwise XOR (exclusive OR) .
*	Multiplication	<=	Less than or equal to
/	Division	>	Greater than

**	Exponentiation	>=	Greater than or equal to.
//	Floor division	!=	Not equal to
<<	bit wise left shift	==	Equal to
>>	right shift	@=	In-place operator (used in some contexts like x @= y, equivalent to x = x @ y).
&	and	:=	Assignment expression
@	decorator or matrix multiplication operator		

Throughout this section, we employ practical examples and exercises to reinforce the concepts covered. Our goal is to ensure that learners not only grasp the basics of Python syntax but also appreciate the language's elegance and readability, setting the stage for their journey into more advanced Python programming for data analysis.

```
# Arithmetic Operators
x = 10
y = 5

addition_result = x + y
print("Addition:", addition_result) # Output: 15

subtraction_result = x - y
print("Subtraction:", subtraction_result) # Output: 5

multiplication_result = x * y
print("Multiplication:", multiplication_result) # Output: 50

division_result = x / y
print("Division:", division_result) # Output: 2.0

exponentiation_result = x ** y
print("Exponentiation:", exponentiation_result) # Output: 100000

floor_division_result = x // y
print("Floor Division:", floor_division_result) # Output: 2

# Bitwise Operators
a = 5 # Binary: 101
b = 3 # Binary: 011

bitwise_and_result = a & b
```

```

print("Bitwise AND:", bitwise_and_result) # Output: 1 (Binary: 001)

bitwise_or_result = a | b
print("Bitwise OR:", bitwise_or_result) # Output: 7 (Binary: 111)

bitwise_xor_result = a ^ b
print("Bitwise XOR:", bitwise_xor_result) # Output: 6 (Binary: 110)

bitwise_left_shift_result = a << 1
print("Bitwise Left Shift:", bitwise_left_shift_result) # Output: 10 (Binary: 1010)

bitwise_right_shift_result = a >> 1
print("Bitwise Right Shift:", bitwise_right_shift_result) # Output: 2 (Binary: 10)

# Comparison Operators
x = 10
y = 5

less_than_or_equal = x <= y
print("Less than or equal to:", less_than_or_equal) # Output: False

greater_than = x > y
print("Greater than:", greater_than) # Output: True

greater_than_or_equal = x >= y
print("Greater than or equal to:", greater_than_or_equal) # Output: True

not_equal = x != y
print("Not equal to:", not_equal) # Output: True

equal_to = x == y
print("Equal to:", equal_to) # Output: False

# Assignment Expression
if (n := len("hello")) > 4:
    print(f"The length of 'hello' is {n}.") # Output: "The length of 'hello' is 5."

# Matrix Multiplication Operator
import numpy as np

matrix_a = np.array([[1, 2], [3, 4]])
matrix_b = np.array([[5, 6], [7, 8]])

```



```
matrix_product = matrix_a @ matrix_b
print("Matrix Product:")
print(matrix_product)
```

3. Control Structures:

Control Structures: Progress to teaching control structures, a critical aspect of Python programming. Begin by explaining conditional statements, including 'if,' 'elif,' and 'else.' These statements enable your learners to create decision-making logic within their programs, a fundamental skill for any data analyst. Additionally, delve into loops, both 'for' and 'while.' Loops allow data analysts to automate repetitive tasks and iterate through data efficiently. Understanding how to use these constructs is essential for managing program flow, a crucial skill when dealing with the intricate workflows encountered in data analysis. Familiarize your learners with real-world examples where control structures are applied to make decisions based on data conditions, emphasizing their practical importance in data-driven scenarios.

```
#Conditional Statements (if-else):
x = 10
if x > 5:
    print("x is greater than 5")
elif x == 5:
    print("x is equal to 5")
else:
    print("x is not greater than 5")
```

while

```
# Example using a while loop
count = 0 # initialize count to 0
while count < 5: # condition: continue to run until condition is false
    print(count) # print count
    count += 1 # increment count by 1
    # this is equivalent to count = count + 1, it will keep adding 1 to count
until count = 5
```

4. Data Structures:

Delve into the realm of data structures, a crucial aspect of Python programming. Begin by introducing fundamental structures such as lists, tuples, dictionaries, and sets. Offer a comprehensive understanding of their distinctive properties, including mutable and immutable aspects, indexing techniques, and how these structures can efficiently store and organize data.

List:

```
# list
fruits = ["apple", "banana", "cherry"]
print(fruits)
# check data type
print(type(fruits))
```

Dictionary

```
#Dictionaries:
person = {"name": "Alice", "age": 25, "city": "New York"}
print(type(person))
```

To enhance comprehension, provide practical insights into how these data structures are used in real-world scenarios. For instance, you can demonstrate how lists are employed to manage collections of data, tuples for representing unchangeable sequences, dictionaries for organizing key-value pairs, and sets for handling unique elements. By illustrating these practical applications, learners can better grasp the relevance and utility of each data structure, making their understanding more tangible and actionable. Additionally, explore various methods for data manipulation within these structures, empowering students to work confidently with data in Python.

5. Functions:

Now, let's delve deeper into Python functions. This section will provide a comprehensive understanding of how to define and call functions, ensuring that learners grasp the nuances of function parameters, return values, and variable scope. Functions play a pivotal role in code organization and facilitate the creation of reusable data analysis processes. They allow you to encapsulate specific operations, making your code modular and efficient. We will explore the versatility of functions, showcasing how they can simplify

complex tasks and enhance the overall readability and maintainability of your data analysis scripts. Additionally, we'll cover advanced topics such as lambda functions and the concept of function decorators, which are valuable tools in the toolkit of a proficient Python data analyst.

```
# Define your own function
def my_square(n):
    return n ** 2
print(my_square(5))
```

6. Modules and Libraries:

In this section; we will delve deeper into the pivotal role that Python's extensive ecosystem of modules and libraries plays in data analysis.

- i. **Python's Standard Library:** Python's Standard Library is a rich collection of modules and packages that come bundled with the language, offering a wide array of functionalities for developers. It encompasses tools for tasks ranging from file handling and regular expressions to network programming and web development. For instance, the `os` module facilitates operating system interactions, while `datetime` manages date and time operations. The `json` module aids in handling JSON data, and `urllib` simplifies URL-related tasks. With `math`, mathematical operations become accessible, and `collections` provides additional data structures. Overall, Python's Standard Library serves as a robust foundation, reducing the need for external dependencies and enabling developers to build versatile applications efficiently.

Some Built-in function

Zip : Combines two or more iterables element-wise into pairs or tuples.

```
names = ["Alice", "Bob", "Charlie"]
scores = [95, 87, 92]

zipped_data = zip(names, scores)
zipped_list = list(zipped_data)
print(zipped_list) #
```

Map: Applies a function to each element of an iterable and returns the results.

```
# use of map without lambda
def square(x):
    return x**2
squared = map(square, numbers)
print(list(squared))

# Using map to square a list of numbers
numbers = [1, 2, 3, 4, 5]
squared = map(lambda x: x**2, numbers) # map(function, iterable)
print(list(squared))
```

filter: Filters elements from an iterable based on a given condition, returning those that meet the condition.

```
numbers = [1, 2, 3, 4, 5, 6, 7, 8, 9]

# Filter even numbers
even_numbers = filter(lambda x: x % 2 == 0, numbers)
even_list = list(even_numbers)
print(even_list) # Output: [2, 4, 6, 8]
```

reduce : Accumulates values in an iterable by applying a binary function repeatedly to reduce them to a single result.

```
from functools import reduce

numbers = [1, 2, 3, 4, 5]
product = reduce(lambda x, y: x * y, numbers)
print(product) # Output: 120
```

- ii. **External Libraries:** Python boasts a vast ecosystem of external libraries that greatly expand its capabilities. These libraries cover a wide range of domains, from data science and machine learning with libraries like NumPy, pandas, and scikit-learn, to web development with Django and Flask. For data visualization, Matplotlib and Seaborn provide powerful tools, while TensorFlow and PyTorch offer advanced deep learning capabilities. Additionally, libraries like Requests simplify HTTP requests, and SQLAlchemy streamlines database interactions. Python's external libraries not only enhance its versatility but also enable developers to tackle diverse projects efficiently by leveraging pre-built solutions and extensive community support.

- iii. **Importing Modules:** Python's modular design encourages code organization and reuse. Developers can import modules, both built-in and third-party, to access pre-written functions, classes, and variables. For instance, by importing the math module, mathematical operations become readily available, and datetime provides tools for managing dates and times. Libraries like NumPy, Pandas, and Matplotlib offer specialized functionality for data manipulation, analysis, and visualization.

```
# import built-in libraries
import math
import random
# importing pandas library, numpy library: Third-party libraries
import pandas as pd
import numpy as np
```

- iv. **Popular Data Analysis Libraries:** Dedicate substantial time to popular data analysis libraries that form the backbone of data manipulation, analysis, and visualization in Python:



Figure 3. Python Popular Data analysis and visualization Libraries.

1. [NumPy](#): we will embark on a journey into the realm of numerical computing with NumPy.

This foundational Python library is crucial for performing efficient numerical operations and

working with arrays. We'll begin by introducing the core data structure of NumPy - the ndarray. You'll learn how to create, manipulate, and perform basic arithmetic operations on arrays. Topics covered will include broadcasting, which allows for element-wise operations on arrays of different shapes, and universal functions (ufuncs) for element-wise operations. You'll also learn about array slicing, indexing, and boolean masking for extracting and manipulating data within arrays efficiently. We will explore techniques for reshaping arrays, stacking and splitting arrays, and aggregating data using functions like sum, mean, and max. Additionally, you'll discover how to handle missing data and perform operations on multidimensional arrays. We will delve into matrix operations, matrix decompositions (e.g., eigenvalues and singular value decomposition), and solving linear equations using NumPy. This section will be invaluable for those involved in scientific computing and machine learning. We will discuss the reasons behind NumPy's performance and introduce techniques for optimizing your code using vectorized operations. You'll also explore NumPy's integration with other libraries and tools, enhancing its utility in diverse applications.

2. **Pandas**: Our journey begins with an introduction to pandas, a powerful library for data manipulation and analysis. We'll learn how to effectively handle data using pandas, covering topics such as importing and exporting data and ensuring seamless data integration into your projects with support for various file formats, including CSV, Excel, and SQL databases. Further, we'll dive into the intricacies of data cleaning and preprocessing, focusing on techniques to ensure data quality and consistency. We'll also explore data filtering, enabling you to extract the relevant information you need. Additionally, we'll delve into data aggregation, grouping, and transformation, where we'll uncover advanced methods for summarizing and reshaping datasets to suit your analytical needs. As a bonus, we'll discuss strategies to speed up your data analysis process using pandas' powerful optimizations. By the end of this session, you'll have a solid grasp of Python's data analysis capabilities, including support for various file formats, merging datasets, grouping, and aggregating data, filtering and

cleaning, reshaping data structures, and techniques to enhance the efficiency of your analytical workflows. You'll be well-equipped to tackle more complex challenges in the days ahead.

3. [Matplotlib](#): we will immerse ourselves in the world of data visualization with Matplotlib, one of Python's most versatile plotting libraries. This session will equip you with the skills needed to create compelling and informative visualizations. We'll start with the basics of creating simple plots, such as line plots, scatter plots, and bar charts, to visualize your data effectively. As we progress, we'll explore advanced techniques for customizing plots, including labels, titles, color palettes, and annotations, ensuring that your visualizations convey meaningful insights. Moreover, we'll delve into the creation of more complex visualizations like histograms, box plots, and heat maps to represent data distributions and relationships. Moving forward, we'll delve into intermediate data visualization topics. We'll explore subplots and figure customization to design multi-panel visualizations that effectively communicate various aspects of your data. You'll also learn to create interactive visualizations that enhance your data exploration experience. Additionally, we'll touch on geographical data visualization with tools like GeoPandas, enabling you to map and analyze spatial data. As a bonus, we'll discuss strategies to optimize your visualizations for publication or presentation.

Matplotlib Function	Description
<code>plt.plot(x, y)</code>	Plot a line graph with data points <code>x</code> and <code>y</code> .
<code>plt.scatter(x, y)</code>	Create a scatter plot with data points <code>x</code> and <code>y</code> .
<code>plt.bar(x, height)</code>	Generate a bar chart with values <code>x</code> and <code>height</code> .
<code>plt.hist(data, bins)</code>	Create a histogram with the given data and bins.
<code>plt.xlabel('xlabel')</code>	Set the label for the x-axis.
<code>plt.ylabel('ylabel')</code>	Set the label for the y-axis.
<code>plt.title('title')</code>	Add a title to the plot.
<code>plt.legend()</code>	Display a legend for labeled elements on the plot.
<code>plt.grid(True)</code>	Enable gridlines on the plot.
<code>plt.savefig('filename.png')</code>	Save the current plot as an image file.

Figure 4. Basic Matplotlib code syntax.

4. **Seaborn:** Highlight Seaborn as an indispensable companion to Matplotlib for creating aesthetically pleasing statistical data visualizations. Showcase Seaborn's capabilities in producing complex plots with minimal code. Explore its features for exploring dataset distributions, relationships, and trends.

Hands-On Practice: Reinforce theoretical knowledge with hands-on exercises and projects that involve using these libraries to solve real-world data analysis challenges. Encourage learners to explore sample datasets and perform tasks like data cleaning, analysis, and visualization using the introduced libraries. By the end of this module, learners will have a solid understanding of how Python's modules and libraries empower data analysts to efficiently process, analyze, and visualize data, making it a crucial skill set in the field of data analysis and beyond.