

# Python for Data Science: A very short introduction for beginners

Jean-Michel Marin

September 19, 2024

## Contents

<b>1</b>	<b>Key differences between Python and R</b>	<b>3</b>
1.1	Primary Purpose . . . . .	3
1.2	Ease of Learning . . . . .	3
1.3	Libraries and Ecosystem . . . . .	3
1.4	Statistical and Data Analysis . . . . .	4
1.5	Data Visualization . . . . .	4
1.6	Machine Learning and AI . . . . .	4
1.7	Community and Support . . . . .	5
1.8	Performance . . . . .	5
1.9	Usage in Industry vs. Academia . . . . .	5
1.10	Integration and Deployment . . . . .	5
1.11	Conclusion . . . . .	5
<b>2</b>	<b>Core Libraries for Data Science in Python</b>	<b>6</b>
<b>3</b>	<b>Numerical Computing with NumPy</b>	<b>6</b>
3.1	Array Creation . . . . .	6
3.2	Array Shape and Reshaping . . . . .	7
3.3	Basic Mathematical Operations . . . . .	7
3.4	Element-wise Operations . . . . .	8
3.5	Matrix Operations . . . . .	8
3.6	Statistical Operations . . . . .	9

3.7	Indexing and Slicing . . . . .	9
3.8	Boolean Indexing . . . . .	9
3.9	Random Numbers . . . . .	10
<b>4</b>	<b>Loops in Python</b>	<b>10</b>
4.1	For Loop . . . . .	10
4.2	While Loop . . . . .	11
<b>5</b>	<b>Functions in Python</b>	<b>11</b>
5.1	Defining and Calling Functions . . . . .	11
5.2	Functions with Return Values . . . . .	11
5.3	Functions with Multiple Parameters . . . . .	12
5.4	Default Parameters in Functions . . . . .	12
<b>6</b>	<b>Data Manipulation with Pandas</b>	<b>12</b>
6.1	Load the CSV File . . . . .	12
6.2	Basic Data Exploration . . . . .	13
6.3	Filtering Data . . . . .	13
6.4	Creating New Columns . . . . .	13
6.5	Grouping and Aggregating Data . . . . .	13
6.6	Sorting Data . . . . .	14
6.7	Handling Missing Data . . . . .	14
6.8	Calculating Descriptive Statistics . . . . .	14
<b>7</b>	<b>Data Visualization with Matplotlib and Seaborn</b>	<b>15</b>
<b>8</b>	<b>Machine Learning with Scikit-learn</b>	<b>15</b>
<b>9</b>	<b>Data Preprocessing</b>	<b>16</b>
<b>10</b>	<b>Deep Learning with TensorFlow or PyTorch</b>	<b>16</b>

Python has become the go-to language for data science due to its simple syntax, extensive libraries, and large community support. This makes it easier for data scientists to clean, manipulate, visualize, and model data with minimal code.

## 1 Key differences between Python and R

Python and R are two of the most popular programming languages used in data science, statistics, and machine learning. While both are powerful and versatile, they have distinct differences in terms of syntax, ecosystem, performance, and usage. Below is a comparison of the two languages:

### 1.1 Primary Purpose

- Python is a general-purpose programming language with a wide range of applications beyond data science, such as web development, automation, and software engineering. It is widely known for its readability and ease of use.
- R was specifically designed for statistical computing and data visualization. It is favored by statisticians, data analysts, and researchers for tasks that involve heavy statistical analysis and visualizing data.

### 1.2 Ease of Learning

- Python is known for its simple and clear syntax, making it easier for beginners to learn. Its readability and similarity to natural language often make Python a preferred choice for those new to programming.
- R has a steeper learning curve compared to Python, especially for users without a programming background. It requires a better understanding of statistical concepts and has a less intuitive syntax compared to Python.

### 1.3 Libraries and Ecosystem

- Python has a vast ecosystem of libraries for data science, such as:
  - Pandas for data manipulation,
  - NumPy for numerical computing,
  - Matplotlib and Seaborn for data visualization,
  - Scikit-learn for machine learning, and
  - TensorFlow and PyTorch for deep learning.

Python is also widely used in web development (with Django, Flask), automation, and artificial intelligence, making it highly versatile.

- R has a rich ecosystem of packages for statistical computing and data visualization, including:
  - ggplot2 for data visualization,
  - dplyr for data manipulation,
  - shiny for building interactive web applications,
  - caret for machine learning.

R's packages are generally more specialized in statistics and data visualization.

## 1.4 Statistical and Data Analysis

- Python is capable of performing statistical analysis, but R has traditionally been the go-to language for complex statistical modeling. Python's statistical libraries (like SciPy and statsmodels) are improving rapidly and are suitable for most tasks, but R still has a deeper range of statistical functions out of the box.
- R is specifically designed for statistical analysis and has a large collection of built-in functions for complex statistical modeling. It is often preferred in academia and research for this reason.

## 1.5 Data Visualization

- Python offers libraries like Matplotlib, Seaborn, and Plotly for data visualization, which are powerful but can sometimes require more effort to create complex plots.
- R excels in data visualization, particularly with its ggplot2 library, which provides an intuitive and highly customizable approach to creating complex and beautiful visualizations with minimal code. R is often favored by data analysts who prioritize visual representation.

## 1.6 Machine Learning and AI

- Python is the dominant language for machine learning and AI, thanks to its comprehensive libraries such as Scikit-learn, TensorFlow, Keras, and PyTorch. Python is preferred in the industry for building machine learning pipelines and deploying models into production environments.
- R has machine learning libraries like caret and randomForest, but it is less widely used for production machine learning systems. R is more focused on exploratory data analysis and statistical modeling.

## 1.7 Community and Support

- Python has an incredibly large and diverse community due to its versatility across different domains (data science, web development, software engineering). As a result, finding tutorials, forums, and open-source projects in Python is relatively easy.
- R also has a strong community, particularly among statisticians and academic researchers. Its community is smaller compared to Python, but highly focused on statistical analysis and data visualization.

## 1.8 Performance

- Python is an interpreted language, which can sometimes make it slower compared to compiled languages. However, Python's libraries (like NumPy, which is implemented in C) ensure that data-intensive tasks can still be performed efficiently.
- R can be slower than Python for large-scale tasks, especially when handling large datasets in memory. However, R has libraries like `data.table` that can speed up operations on large datasets.

## 1.9 Usage in Industry vs. Academia

- Python is widely used in both industry and academia. In industry, Python dominates in fields such as machine learning, AI, data engineering, and software development. Python's ability to integrate into production environments makes it ideal for deployment and scalability.
- R is more commonly used in academia and research, particularly for statistical analysis. It's also widely used in industries like bioinformatics, healthcare, and economics where statistical analysis is the primary task.

## 1.10 Integration and Deployment

- Python integrates easily with other systems and platforms, making it suitable for deploying machine learning models, building APIs, and creating end-to-end data pipelines. Python is often preferred for full-stack data science work (from data analysis to model deployment).
- R is more focused on data analysis and reporting. While it has tools like Shiny for creating interactive dashboards, it's not as commonly used for deployment into production systems compared to Python.

## 1.11 Conclusion

- Choose Python if you are looking for a general-purpose language that excels in machine learning, AI, and production environments. It's great for tasks that involve data engineering, automation, and deploying models.

- Choose R if you are focused on statistical analysis, data visualization, or academic research. R's ecosystem is built around statistical methods, making it easier for statisticians and data analysts to work with complex datasets.

## 2 Core Libraries for Data Science in Python

To get started with Python for data science, you'll need to work with some core libraries:

- **NumPy**: Supports large multi-dimensional arrays and matrices, along with a collection of mathematical functions to operate on these arrays.
- **Pandas**: Built on top of NumPy, it's the key library for data manipulation and analysis.
- **Matplotlib and Seaborn**: These are visualization libraries that make it easy to create a wide range of plots.
- **Scikit-learn**: A comprehensive machine learning library.
- **TensorFlow and PyTorch**: Libraries for deep learning, although they are more advanced.

## 3 Numerical Computing with NumPy

NumPy is essential for numerical operations, particularly for working with large arrays and matrices.

### 3.1 Array Creation

```
def square(number): return number**2
# Call the function and print the result
result = square(5)
print(f'The square of 5 is {result}')
```

---

```
import numpy as np
# Create a 1D array
arr_1d = np.array([1, 2, 3, 4, 5])
print(arr_1d)
# Create a 2D array (matrix)
arr_2d = np.array([[1, 2, 3], [4, 5, 6]])
print(arr_2d)
# Create an array of zeros
zeros = np.zeros((3, 3))
print(zeros)
# Create an array of ones
ones = np.ones((2, 4))
```

```
print(ones)
# Create an array with a range of values
range_array = np.arange(0, 10, 2)
print(range_array)
# Create an array with equally spaced values (linspace)
linspace_array = np.linspace(0, 1, 5)
print(linspace_array)
# Create a random array
random_array = np.random.rand(3, 3)
print(random_array)
# Create an identity matrix
identity_matrix = np.eye(3)
print(identity_matrix)
```

---

### 3.2 Array Shape and Reshaping

---

```
# Check the shape of an array
print(arr_2d.shape) # Output: (2, 3)
# Reshape an array
reshaped_array = arr_1d.reshape((5, 1))
print(reshaped_array)
# Flatten a 2D array into a 1D array
flattened_array = arr_2d.flatten()
print(flattened_array)
```

---

### 3.3 Basic Mathematical Operations

---

```
arr = np.array([10, 20, 30, 40])
# Add 5 to each element
arr_add = arr + 5
print(arr_add)
# Subtract 5 from each element
arr_sub = arr - 5
print(arr_sub)
# Multiply each element by 2
arr_mul = arr * 2
print(arr_mul)
# Divide each element by 2
arr_div = arr / 2
print(arr_div)
# Exponentiation (square each element)
arr_square = np.square(arr)
print(arr_square)
```

```
# Square root of each element
arr_sqrt = np.sqrt(arr)
print(arr_sqrt)
```

---

### 3.4 Element-wise Operations

---

```
arr_a = np.array([1, 2, 3])
arr_b = np.array([4, 5, 6])
# Element-wise addition
result_add = arr_a + arr_b
print(result_add)
# Element-wise multiplication
result_mul = arr_a * arr_b
print(result_mul)
# Element-wise comparison
result_comp = arr_a > arr_b
print(result_comp)
# Sum of all elements in an array
total_sum = np.sum(arr_a)
print(total_sum)
# Sum along an axis (rows or columns)
matrix = np.array([[1, 2, 3], [4, 5, 6]])
sum_rows = np.sum(matrix, axis=0) # Sum along columns
sum_cols = np.sum(matrix, axis=1) # Sum along rows
print(sum_rows, sum_cols)
```

---

### 3.5 Matrix Operations

---

```
matrix_a = np.array([[1, 2], [3, 4]])
matrix_b = np.array([[5, 6], [7, 8]])
# Matrix multiplication
matrix_product = np.dot(matrix_a, matrix_b)
print(matrix_product)
# Transpose of a matrix
transpose_a = matrix_a.T
print(transpose_a)
# Inverse of a matrix
inverse_a = np.linalg.inv(matrix_a)
print(inverse_a)
# Determinant of a matrix
det_a = np.linalg.det(matrix_a)
print(det_a)
```

---



### 3.6 Statistical Operations

---

```
data = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9])
# Mean
mean = np.mean(data)
print(mean)
# Median
median = np.median(data)
print(median)
# Standard deviation
std_dev = np.std(data)
print(std_dev)
# Variance
variance = np.var(data)
print(variance)
# Min and Max
min_val = np.min(data)
max_val = np.max(data)
print(min_val, max_val)
# Percentile
percentile_50 = np.percentile(data, 50)
print(percentile_50)
```

---

### 3.7 Indexing and Slicing

---

```
arr = np.array([10, 20, 30, 40, 50])
# Access the first element
print(arr[0])
# Access the last element
print(arr[-1])
# Slice the array (from index 1 to 3)
print(arr[1:4])
# Reverse the array
print(arr[::-1])
# Modify elements
arr[1:3] = [25, 35]
print(arr)
```

---

### 3.8 Boolean Indexing

---

```
arr = np.array([10, 15, 20, 25, 30])
# Create a boolean array where elements are greater than 20
bool_idx = arr > 20
```

```
print(bool_idx)
# Filter elements based on the boolean array
filtered_arr = arr[bool_idx]
print(filtered_arr)
# Set elements greater than 20 to 100
arr[arr > 20] = 100
print(arr)
```

---

### 3.9 Random Numbers

---

```
# Generate random numbers between 0 and 1
random_nums = np.random.rand(5)
print(random_nums)
# Generate random integers between a range
random_ints = np.random.randint(1, 10, size=5)
print(random_ints)
# Generate random numbers from a normal distribution
random_normal = np.random.randn(5)
print(random_normal)
# Seed the random number generator for reproducibility
np.random.seed(42)
random_seeded = np.random.rand(3)
print(random_seeded)
```

---

## 4 Loops in Python

### 4.1 For Loop

A for loop is used to iterate over a sequence (list, tuple, string, or range).

---

```
# Iterate over a list
my_list = [1, 2, 3, 4, 5]
for item in my_list:
    print(item)
# Iterate over a range of numbers
for i in range(5):
    print(f'Value: {i}')
```

---

In the above example, the for loop iterates over the elements in `my_list` and prints each value. The second example shows how to use the `range()` function to loop through numbers 0 to 4.

## 4.2 While Loop

A `while` loop runs as long as the condition is true.

---

```
# Initialize a variable
counter = 0
# Run the loop while the condition is true
while counter < 5:
    print(f'Counter is: {counter}')
    counter += 1
```

---

In this example, the loop continues to run while the variable `counter` is less than 5. On each iteration, `counter` is incremented by 1 until the condition is no longer true.

## 5 Functions in Python

### 5.1 Defining and Calling Functions

Functions allow you to define reusable blocks of code. You define a function using the `def` keyword.

---

```
# Define a simple function
def greet(name):
    print(f'Hello, {name}!')
# Call the function
greet('Alice')
greet('Bob')
```

---

In the above example, `greet()` is a function that takes a parameter `name` and prints a greeting message. The function is called with different arguments.

### 5.2 Functions with Return Values

You can also define functions that return values.

---

```
def square(number): return number**2
# Call the function and print the result
result = square(5)
print(f'The square of 5 is {result}')
```

---

Here, the `square()` function returns the square of the input value. The result of the function call is stored in the variable `result` and printed.

## 5.3 Functions with Multiple Parameters

Functions can take multiple parameters.

---

```
# Define a function with multiple parameters
def add_numbers(a, b):
    return a + b
# Call the function
sum_result = add_numbers(10, 20)
print(f'The sum of 10 and 20 is {sum_result}')
```

---

In this example, the `add_numbers()` function takes two arguments, adds them, and returns the result.

## 5.4 Default Parameters in Functions

You can specify default values for function parameters.

---

```
# Define a function with a default parameter
def greet(name="Guest"):
    print(f'Hello, {name}! ')
# Call the function without passing an argument
greet()
# Call the function with an argument
greet('Alice')
```

---

In the above code, the function `greet()` has a default parameter value of "Guest". If no argument is passed when calling the function, the default value is used.

# 6 Data Manipulation with Pandas

Pandas is used for handling and analyzing structured data. It provides two main data structures: Series (1D data) and DataFrame (2D data, similar to an Excel sheet or SQL table).

## 6.1 Load the CSV File

---

```
import pandas as pd
# Load the csv file
df = pd.read_csv('data.csv')
# Display the DataFrame
print(df)
```

---

## 6.2 Basic Data Exploration

---

```
# Get the first few rows of the dataset
print(df.head())
# Get the data types of each column
print(df.dtypes)
# Get a summary of the numerical columns
print(df.describe())
# Get the number of missing values for each column
print(df.isnull().sum())
```

---

## 6.3 Filtering Data

---

```
# Filter rows where Column1 > 200
filtered_data = df[df['Column1'] > 200]
print(filtered_data)
# Filter rows where CategoricalColumn is 'A'
category_A = df[df['CategoricalColumn'] == 'A']
print(category_A)
```

---

## 6.4 Creating New Columns

---

```
# Create a new column that is the product of
# Column1 and Column2
df['Column3'] = df['Column1'] * df['Column2']
print(df)
# Create a new column that categorizes based on Column1 values
df['Column1_Category'] = df['Column1'].
apply(lambda x: 'High' if x > 250 else 'Low')
print(df)
```

---

## 6.5 Grouping and Aggregating Data

---

```
# Group by CategoricalColumn and calculate
# the mean of Column1 for each group
grouped = df.groupby('CategoricalColumn')['Column1'].mean()
print(grouped)
# Group by CategoricalColumn and
# get multiple aggregations for Column1
aggregated = df.groupby('CategoricalColumn')['Column1'].
agg(['mean', 'sum', 'max', 'min'])
print(aggregated)
```

---

## 6.6 Sorting Data

---

```
# Sort the DataFrame by Column1 in descending order
sorted_df = df.sort_values(by='Column1', ascending=False)
print(sorted_df)
# Sort by multiple columns: first by CategoricalColumn,
# then by Column1
sorted_multi = df.sort_values(by=['CategoricalColumn', 'Column1'],
    ascending=[True, False])
print(sorted_multi)
```

---

## 6.7 Handling Missing Data

---

```
# Replace missing values in Column1 with the mean of the column
df['Column1'].fillna(df['Column1'].mean(), inplace=True)
# Drop rows with missing values
df.dropna(inplace=True)
```

---

## 6.8 Calculating Descriptive Statistics

---

```
# Create a sample DataFrame
data = {
    'Age': [25, 30, 35, 40, 45],
    'Height': [165, 170, 175, 180, 185],
    'Weight': [65, 70, 75, 80, 85]
}
df = pd.DataFrame(data)
# Display the DataFrame
print(df)

# Calculate the mean of each column
mean_values = df.mean()
print("Mean Values:")
print(mean_values)

# Calculate the median of each column
median_values = df.median()
print("Median Values:")
print(median_values)

# Calculate the standard deviation of each column
std_dev_values = df.std()
print("Standard Deviation Values:")
```

```
print(std_dev_values)

# Calculate the variance of each column
variance_values = df.var()
print("Variance Values:")
print(variance_values)

# Calculate the correlation between columns
correlation_matrix = df.corr()
print("Correlation Matrix:")
print(correlation_matrix)

# Generate a summary of descriptive statistics
summary_stats = df.describe()
print("Summary Statistics:")
print(summary_stats)
```

---

## 7 Data Visualization with Matplotlib and Seaborn

Visualization is key in data science to understand trends, patterns, and correlations. Matplotlib and Seaborn are the two most commonly used libraries for data visualization.

---

```
import matplotlib.pyplot as plt
import seaborn as sns
plt.plot([1, 2, 3, 4], [1, 4, 9, 16])
plt.xlabel('X Axis')
plt.ylabel('Y Axis')
plt.title('Simple Line Plot')
plt.show()
sns.histplot(df['Column1'], kde=True)
plt.title('Distribution of Column1')
plt.show()
corr_matrix = df.corr()
sns.heatmap(corr_matrix, annot=True)
plt.title('Correlation Matrix Heatmap')
plt.show()
```

---

## 8 Machine Learning with Scikit-learn

Scikit-learn is the go-to library for machine learning in Python. It provides easy-to-use tools for data pre-processing, building models, and evaluating them.

---

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
X = df[['Column1', 'Column2']]
y = df['Target']
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)
model = LinearRegression()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
print(f'Mean Squared Error: {mse}')
```

---

## 9 Data Preprocessing

Before building a machine learning model, data preprocessing is essential. Scikit-learn provides a variety of preprocessing tools:

---

```
from sklearn.preprocessing import StandardScaler, LabelEncoder
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
label_encoder = LabelEncoder()
df['CategoricalColumn'] = label_encoder.
fit_transform(df['CategoricalColumn'])
```

---

## 10 Deep Learning with TensorFlow or PyTorch

For deep learning tasks, such as image recognition or natural language processing, you'll often use TensorFlow or PyTorch. These libraries allow you to build neural networks and train models on large datasets. Here's an example with TensorFlow:

---

```
import tensorflow as tf
model = tf.keras.Sequential([
    tf.keras.layers.Dense(64, activation='relu',
input_shape=(X_train.shape[1],)),
    tf.keras.layers.Dense(32, activation='relu'),
    tf.keras.layers.Dense(1)
])
model.compile(optimizer='adam', loss='mean_squared_error')
model.fit(X_train, y_train, epochs=10, batch_size=32)
loss = model.evaluate(X_test, y_test)
print(f'Loss: {loss}')
```

---