

# Session 1. Basics of working with Python

2026-02-13

---

## Highlights:

This is my mini-reflection. Paragraphs must be indented.  
It can contain multiple paragraphs.

---

“Be willing to be a beginner every single morning.”

— Meister Eckhart

## Session outline

- Why Python?
- Installing the software: Python and RStudio (with `reticulate`)
- Packages (libraries)
- Installing packages: PyPI, Conda, and other sources
- Getting help
- Creating a project
- Directory structure
- Creating new files: types of files
- Literate programming
- Data objects/classes and basic operations

## Introduction

NOTE: This is an Quarto Markdown document. This type of document is a plain text file that can recognize chunks of code. When you execute code within the document, the results are displayed beneath. Quarto Markdown files are *computational notebooks* which implement a coding philosophy called *literate programming*. Literate computing emphasizes the use of natural language to communicate with humans and chunks of code to communicate with the computer. By making the main audience other humans, this style of coding flips around the usual way in which code is written (computer is main audience, humans come second). This helps to make learning how to code more intuitive and accessible.

Let us first explain the anatomy of an Quarto Markdown document.

At the top of the document is a *header* (called a YAML header) contained between two sets of three dashes ---. The header includes metainformation about the document, including possibly the title of the document, the name(s) of the author(s), and how to process the document (e.g., what type of output to produce). After the header comes the *body* of the document, that is, the main contents of the document (text and code).

The header is quite an important component of the document, but we will for the moment take it as given and concentrate on the body of the document. This file, for example, introduces some basic instructions to work with Python and Quarto Markdown.

Chunks of code are key to writing computational notebooks. Whenever you see a chunk of code in an Quarto Markdown document as follows, you can run it (by clicking the green ‘play’ icon on the top right corner of the code window) to see the results. Try it below!

```
print("Hello, Workshop")
```

Hello, Workshop

You can see that the function `print()` displays the argument on the screen.

Notice the way the document speaks to you in natural language, and to the computer in Python. The computer, instead of being the focus of the document, is an assistant to illustrate concepts (like, what is a chunk of code?).

Each of the documents in this workshop (the *Sessions*) is like a set of lecture notes. You can *knit* the document to produce a PDF file to study. The documents, on the other hand, are potentially much more than simple lecture notes, and they can in fact become the foundation for your own experiments and annotations. As you read and study, you can ‘customize’ the notes based on your developing understanding of the subject matter and/or to complement the document with other examples and information. To make the readings uniquely yours, you can type directly into the document (the original template is still available from the package, and if you need a fresh start, you can always create a new file with it).

In addition, you can use the following style to create a *textbox*:

**NOTICE:**

This is an annotation. Here I can write my thoughts as I study, or I can add useful links or other information to help me learn.

To create a new paragraph, I need to type two blanks after the last one.

A textbox allows you to highlight important information. Try creating your own textbox next.

By now you will already have noticed a few conventions for writing in R Markdown:

- A chunk of Python code begins with three backticks (```) and the letter ‘python’ between curly brackets; a chunk of code concludes with three backticks.
- Hashtags (#) indicate headers: one hashtag is a primary header, two hashtags is a secondary header, three a tertiary header, etc.
- Asterisks are used for changing the font to *Italics*, **Bold**, and ***Italics+Bold***.
- Superscrit can be added by using ^, such as in superscript<sup>2</sup>
- Subscript can be added by using ~, such as in superscript<sub>2</sub>
- Dashes are used to create unnumbered lists.
- Hyperlinks can be introduced by using square brackets followed by the url in brackets.
- Greater-than (>) is used for block-quotes.

There are some other typing conventions that you can find in this useful cheatsheet.

## RStudio Window

If you are reading this, you must already have learned how to create a new project and documents using the templates in the workshop package. We can now proceed to discuss some basic concepts regarding data types and operations.

## Preliminaries

We will load all *libraries* we plan to use in this session. A library (or package) is a collection of modules containing pre-written code, functions, and classes that developers can use to perform common tasks without writing code from scratch. These are shareable code units that extend base Python’s capabilities. When you install Python, it includes built-in standard libraries for mathematical operations, date/time handling, system interactions, and more.

However, sometimes we need to perform analyses for which no built-in libraries exist. In these cases, we must install libraries from external sources hosting package repositories or “indexes.” The two most common are PyPI(Python Package Index) and Conda repositories, but packages can also come from GitHub repositories and other locations.

Conda’s advantage is its dependency tracking—it manages how libraries depend on one another, reducing compatibility issues.

For PyPI packages, use `pip install package_name` in the terminal. For Conda, use `conda install package_name`.

Two libraries will be extensively used in this workshop: `matplotlib` and `pandas`. `Matplotlib` is Python’s most popular data visualization library, enabling the creation of 2D plots. `Pandas` provides robust data structures and functions for working with tabular data, primarily Series and DataFrames.

If you haven’t already, install these packages with `pip install pandas matplotlib` or `conda install pandas matplotlib`. Once installed, you must load them into memory using `import package_name`:

```
import pandas as pd
```

## Language Semantics

Python uses indentation (spaces or tabs) to structure code, unlike languages like R and Java, which use curly braces ({}).

Text following the # symbol is treated as a comment and ignored by the Python interpreter:

```
data # Any comment  
# Another comment
```

As in other languages, Python functions are reusable code blocks that perform specific actions when called. We call them by name followed by parentheses (()). Parameters can be included within the parentheses. For example, a function `f` with parameters `a`, `b`, and `c` is called as:

```
f(a, b, c)
```

In Python, everything is an *object*. Every data type is a Python object with a specific type, internal data, and associated functions. Functions within objects are called *methods*, providing access to the object’s content. For example, if we have an object `data` and want to access information via its method:

```
data.method()
```

Objects also have *attributes*: values representing properties or characteristics. Like methods, attributes are accessed using dot notation:

```
data.attribute
```

Python uses = to assign values to variables. Variables are symbolic names that reference objects or values stored in memory:

```
var = 10  
var
```

10

In the chunk above, we assigned the value 10 to the variable `var`. This assignment allows us to retrieve the value later.

## Data types in Python

After executing the previous chunk, you'll notice a variable named `var` appears in your Environment (upper-right pane tab). This variable holds a single value (10) and belongs to a group of data types known as scalars, which store single values. Standard built-in scalar types are:

Table 1: Standard scalars in Python.

Type	Description
<code>int</code>	Integer (whole numbers)
<code>float</code>	Floating-point numbers, to represent real numbers with a decimal point.
<code>bool</code>	A boolean value. True or False.
<code>bytes</code>	Pure Bytes ASCII
<code>None</code>	Null value
<code>str</code>	String type, to represent text data.

You can check a variable's data type using the `type()` function:

```
type(var)  
  
<class 'int'>
```

Next, we'll explore these scalar data types in more detail.

### Scalars

#### Numeric (`int` and `float`)

The main numeric types in Python are `int` and `float`:

```
1  
  
type(1)  
  
<class 'int'>  
  
3.141593  
  
3.141593  
  
type(3.141593)  
  
<class 'float'>
```

#### Boolean

Boolean values are `True` and `False`:

```
False  
  
False
```

```
type(False)
```

```
<class 'bool'>
```

You can perform Boolean comparisons using `and`, `or`, and comparison operators (`==` for equal, `!=` for not equal, `>` for greater than, `>=` for greater than or equal, `<` for less than, `<=` for less than or equal):

```
4 > 3
```

```
True
```

```
2 == 3
```

```
False
```

```
5 == 2
```

```
False
```

```
4 >= 3
```

```
True
```

### *Strings*

Strings represent text and are created using quotation marks (either `"` or `'`):

```
"This is a string"
```

```
'This is a string'
```

Numbers within quotation marks are stored as strings:

```
type("5")
```

```
<class 'str'>
```

Multi-line strings use triple quotes:

```
"""
This
is
an
example
of
a
string
with
multiple
lines
"""
```

```
'\nThis \nis \nan \nexample\nof \na \nstring\nwith\nmultiple \nlines\n'
```

The `\n` character indicates line breaks in text.

*None*

`None` represents null value in Python:

`None`

### *Built-in data structure*

Now, consider storing multiple values in a variable. For this, you need data structures beyond scalars. We'll explore Python's three built-in data structures: `tuples`, `lists`, and `dictionaries`.

#### *Tuple*

A tuple is a finite, ordered collection of items. Create a tuple by separating items with commas:

```
a_tuple = 100, 150, 45  
a_tuple
```

(100, 150, 45)

You can create nested tuples:

```
nestedtuple = (18, 19, 20), (5, 4)  
nestedtuple
```

((18, 19, 20), (5, 4))

You can also create tuples using the `tuple()` function:

```
tuple([2, 5, 10])
```

(2, 5, 10)

Tuples can store different scalar types:

```
tuple([2, 5, 10, "b"])
```

(2, 5, 10, 'b')

You can create tuples from any iterable sequence:

```
tuple("example")
```

('e', 'x', 'a', 'm', 'p', 'l', 'e')

Tuple elements have fixed positions and are immutable after creation. This structure is useful for storing related data of different types efficiently, such as coordinates (x, y) or user information (name, age, ID).

#### *List*

Lists are similar to vectors in other languages but can store different data types within the same list.

Create a list using square brackets:

```
b_list = [45, 50, 55, 60]
b_list
```

```
[45, 50, 55, 60]
```

Create a list using square brackets:

```
a_list = list(a_tuple)
a_list
```

```
[100, 150, 45]
```

Unlike tuples, lists are mutable. You can append items:

```
a_list.append(200)
a_list
```

```
[100, 150, 45, 200]
```

Insert items at specific positions:

```
b_list.insert(0, 40) # Add 40 at position 0 (first position)
b_list
```

```
[40, 45, 50, 55, 60]
```

Remove items by index using `pop()`:

```
b_list.pop(0) # Removing the value 40 (index 0) from the list
```

```
40
```

```
b_list
```

```
[45, 50, 55, 60]
```

Remove specific items by value using `remove()`:

```
b_list.remove(50)
b_list
```

```
[45, 55, 60]
```

### Dictionaries

A dictionary (`dict`) is a collection of key-value pairs with flexible size, where keys and values are Python objects. Create a dictionary using curly braces with key-value pairs separated by colons:

```
fruits_dict = {'orange': 10, "banana": [5, 10, 20]}
fruits_dict
```

```
{'orange': 10, 'banana': [5, 10, 20]}
```

Add new key-value pairs:

```
fruits_dict["lemon"] = [40, 60, 80]
fruits_dict

{'orange': 10, 'banana': [5, 10, 20], 'lemon': [40, 60, 80]}
```

Access values by key:

```
fruits_dict["banana"]
```

```
[5, 10, 20]
```

Remove a key-value pair using pop():

```
fruits_dict.pop("banana")
```

```
[5, 10, 20]
```

```
fruits_dict
```

```
{'orange': 10, 'lemon': [40, 60, 80]}
```

List all keys in a dict:

```
fruits_dict.keys()
```

```
dict_keys(['orange', 'lemon'])
```

List all values in a dict:

```
fruits_dict.values()
```

```
dict_values([10, [40, 60, 80]])
```

## *Indexing*

Indexing in Python accesses individual items in sequences (strings, lists, tuples) by their position. Python uses zero-based indexing: the first element is at index 0, the second at 1, and so on. Use square brackets [] after the sequence name.

Two indexing approaches:

- Positive indexing: Counts from the beginning (starting at 0).
- Negative indexing: Counts from the end (starting at -1 for the last item).

Retrieve the third character of a string:

```
"This is a string"[2]
```

```
'i'
```

Indexing also works for tuples:

```
a_list
```

```
[100, 150, 45, 200]
```

```
a_list[0] # the indexes starts on 0!
```

```
100
```

```
a_list[1]
```

```
150
```

```
a_list[2]
```

```
45
```

And for lists:

```
b_list[-1]
```

```
60
```

For dictionaries with multiple values per key, access specific values using indexing within the values:

```
fruits_dict['lemon'][0] # for the values value
```

```
40
```

```
fruits_dict['lemon'][0:2]
```

```
[40, 60]
```

## Series and DataFrames

Other useful data structures in Python are Series and DataFrames, included in the `pandas` package. A Series is an array-like object containing a sequence of values and an associated array of labels called an index:

```
materials = pd.Series(["pens", "pencils", "paper", "notebook"])
materials
```

```
0      pens
1    pencils
2      paper
3   notebook
dtype: object
```

We create a Series by calling `pd.Series()` with the elements inside parentheses. The materials Series contains school supplies with indexes on the left. Since we didn't specify index values, pandas created integer indexes starting from 0.

Access Series values by index:

```
materials[0] # For one value
```

```
'pens'
```

```
materials[0:2] # For more than one value, on a sequence
```

```
0      pens
1    pencils
dtype: object
```

```
materials[[1, 0, 2]] # For specific values and order
```

```
1    pencils
0      pens
2    paper
dtype: object
```

Create a Series from a dictionary-keys become indexes, values become Series values:

```
new_series = pd.Series(fruits_dict)
new_series
```

```
orange          10
lemon    [40, 60, 80]
dtype: object
```

For numerical Series, use NumPy for mathematical operations. NumPy is a Python library for numerical computing:

```
import numpy as np

numerical_series = pd.Series([0.30, 0.33, 0.36, 0.39])
np.exp(numerical_series) # For exponential
```

```
0    1.349859
1    1.390968
2    1.433329
3    1.476981
dtype: float64
```

```
numerical_series * 2 # For multiplication
```

```
0    0.60
1    0.66
2    0.72
3    0.78
dtype: float64
```

Note that indexes are preserved.

Logical operations work on numerical Series:

```
numerical_series > 0.35
```

```
0    False
1    False
2    True
3    True
dtype: bool
```

And also for non-numerical Series:

```
'calculator' in materials
```

```
False
```

A DataFrame is a rectangular, tabular data structure consisting of rows and columns (essentially a collection of Series). DataFrames store data in digital format, similar to spreadsheet tables. DataFrames can handle large amounts of information (billions of items, depending on computer memory). Data can be numeric, string, Boolean, etc. Each cell has an address identified by its row and column. DataFrames have indexes for rows and columns.

To illustrate a DataFrame, we will create a dictionary with the vectors for regions in Northern Italy. These data come from a survey of Ph.D. students expected to finish their studies between 2008 and 2014 (see Muscio and Ramaciotti, 2018). The first vector includes the names of the regions; the second vector contains the number of respondents to the survey; third is the number of active spinoffs founded by the students surveyed; next is the number of employees in those spinoffs; and lastly, there is a vector with the total number of active spinoffs in the regions in the period 2005-2006:

```
italy_data = {
    'Region': ["Emilia-Romagna",
               "Friuli-Venezia Giulia",
               "Liguria",
               "Lombardia",
               "Piemonte",
               "Trentino-Alto Adige",
               "Veneto"],
    'PhD_Students': [82, 26, 21, 115, 50, 9, 73],
    'Active_Spinoff': [3, 4, 3, 5, 4, 0, 2],
    'Employees': [13, 7, 8, 3, 16, 0, 1],
    'Spinoff0506': [393, 68, 21, 313, 322, 18, 306]
}

df = pd.DataFrame(italy_data)
df
```

	Region	PhD_Students	Active_Spinoff	Employees	Spinoff0506
0	Emilia-Romagna	82	3	13	393
1	Friuli-Venezia Giulia	26	4	7	68
2	Liguria	21	3	8	21
3	Lombardia	115	5	3	313
4	Piemonte	50	4	16	322
5	Trentino-Alto Adige	9	0	0	18
6	Veneto	73	2	1	306

Different from the previous scalars and data structures, in RStudio, you can click on a DataFrame in the Environment (upper right pane tab) to visualize the data as tabular data in a new tab. You can sort and filter values like in a spreadsheet.

DataFrames are rectangular: all columns must have the same length because each row represents an object of interest (here, a region). If information is missing for some rows, we must either exclude those rows or mark missing values appropriately.

There's one exception: when creating a DataFrame from a dictionary with different value lengths, pandas creates rows equal to the maximum length and fills missing values appropriately. For example, using fruits\_dict (with 2 keys: 'orange' has 1 value, 'lemon' has 3 values):

```
fruits_dict
{'orange': 10, 'lemon': [40, 60, 80]}
```

```
pd.DataFrame(fruits_dict)
```

```
orange  lemon
0      10     40
1      10     60
2      10     80
```

For more than two columns, the table must be rectangular. *You cannot create a DataFrame from a dictionary with keys of different lengths.*

Returning to our municipalities DataFrame, view the first five rows with the method `head()`:

```
df.head()
```

```
          Region  PhD_Students  Active_Spinoff  Employees  Spinoff0506
0  Emilia-Romagna           82            3         13        393
1  Friuli-Venezia Giulia        26            4          7        68
2       Liguria           21            3          8        21
3    Lombardia           115            5          3       313
4     Piemonte           50            4         16        322
```

Get general information with `.info()`, which shows index type, column names, data types, non-null counts, and memory usage:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7 entries, 0 to 6
Data columns (total 5 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Region            7 non-null     object 
 1   PhD_Students       7 non-null     int64  
 2   Active_Spinoff     7 non-null     int64  
 3   Employees          7 non-null     int64  
 4   Spinoff0506        7 non-null     int64  
dtypes: int64(4), object(1)
memory usage: 408.0+ bytes
```

You can obtain a statistical summary by applying the method `.describe()`. For numeric columns, the describe output contains the count (number of non-null observations), mean, standard deviation, minimum value, percentile values and maximum value.

```
df.describe()
```

```
          PhD_Students  Active_Spinoff  Employees  Spinoff0506
count      7.000000     7.000000    7.000000    7.000000
mean      53.714286     3.000000   6.857143  205.857143
std       38.252295    1.632993   6.039552 162.525309
min       9.000000     0.000000   0.000000  18.000000
25%      23.500000     2.500000   2.000000  44.500000
50%      50.000000     3.000000   7.000000 306.000000
75%      77.500000     4.000000  10.500000 317.500000
max      115.000000    5.000000  16.000000 393.000000
```

By default, `describe()` analyzes only numeric columns. To include string/categorical columns, use `include = 'object'` or `include = 'all'`. For object data, the summary includes count, unique values, most frequent value (top), and its frequency:

```
df.describe(include='all')
```

	Region	PhD_Students	Active_Spinoff	Employees	Spinoff0506
count	7	7.000000	7.000000	7.000000	7.000000
unique	7	NaN	NaN	NaN	NaN
top	Emilia-Romagna	NaN	NaN	NaN	NaN
freq	1	NaN	NaN	NaN	NaN
mean	NaN	53.714286	3.000000	6.857143	205.857143
std	NaN	38.252295	1.632993	6.039552	162.525309
min	NaN	9.000000	0.000000	0.000000	18.000000
25%	NaN	23.500000	2.500000	2.000000	44.500000
50%	NaN	50.000000	3.000000	7.000000	306.000000
75%	NaN	77.500000	4.000000	10.500000	317.500000
max	NaN	115.000000	5.000000	16.000000	393.000000

Access DataFrame columns by name within square brackets:

```
df['Region']
```

0	Emilia-Romagna
1	Friuli-Venezia Giulia
2	Liguria
3	Lombardia
4	Piemonte
5	Trentino-Alto Adige
6	Veneto

Name: Region, dtype: object

Or using dot notation (attribute access, not recommended if column names contain spaces or conflict with DataFrame methods):

```
df.Region
```

0	Emilia-Romagna
1	Friuli-Venezia Giulia
2	Liguria
3	Lombardia
4	Piemonte
5	Trentino-Alto Adige
6	Veneto

Name: Region, dtype: object

And for more than one column:

```
df[['Region', 'Spinoff0506']]
```

	Region	Spinoff0506
0	Emilia-Romagna	393
1	Friuli-Venezia Giulia	68
2	Liguria	21

```

3           Lombardia      313
4           Piemonte      322
5   Trentino-Alto Adige      18
6           Veneto       306

```

To access a specific cell:

```

df['Region'][5] # For the row of index 5

'Trentino-Alto Adige'

```

To access an entire row, use .loc:

```

df.loc[5] # For the row of index 5

Region          Trentino-Alto Adige
PhD_Students            9
Active_Spinoff          0
Employees                0
Spinoff0506             18
Name: 5, dtype: object

```

We can create a column or update values in a column that already exists by doing:

```

df["Test"] = 10
df

      Region  PhD_Students  ...  Spinoff0506  Test
0   Emilia-Romagna      82  ...      393     10
1 Friuli-Venezia Giulia      26  ...       68     10
2          Liguria      21  ...       21     10
3          Lombardia     115  ...      313     10
4          Piemonte      50  ...      322     10
5   Trentino-Alto Adige      9  ...       18     10
6          Veneto       73  ...      306     10

```

[7 rows x 6 columns]

In the case above, a new column “Test” is created with all values set to 10.

Add new rows using concat(). In this example, the new dictionary only has values for the “Region” column, so other columns are filled with NaN:

```

new_rows = [{"Region": "Test"}]

df = pd.concat([df, pd.DataFrame(new_rows)], ignore_index=True)

df

```

```

      Region  PhD_Students  ...  Spinoff0506  Test
0   Emilia-Romagna      82.0  ...      393.0  10.0
1 Friuli-Venezia Giulia      26.0  ...       68.0  10.0
2          Liguria      21.0  ...       21.0  10.0
3          Lombardia     115.0  ...      313.0  10.0
4          Piemonte      50.0  ...      322.0  10.0
5   Trentino-Alto Adige      9.0  ...       18.0  10.0

```

```
6          Veneto      73.0   ...
7          Test       NaN   ...

```

[8 rows x 6 columns]

Remove a column with `drop()` and `axis = 1`:

```
df = df.drop('Test', axis = 1)
df
```

```
          Region  PhD_Students  Active_Spinoff  Employees  Spinoff0506
0    Emilia-Romagna        82.0          3.0       13.0      393.0
1  Friuli-Venezia Giulia     26.0          4.0       7.0       68.0
2        Liguria           21.0          3.0       8.0       21.0
3      Lombardia          115.0          5.0       3.0      313.0
4      Piemonte           50.0          4.0      16.0      322.0
5  Trentino-Alto Adige      9.0          0.0       0.0       18.0
6        Veneto           73.0          2.0       1.0      306.0
7          Test            NaN         NaN       NaN       NaN
```

Remove a row with `axis = 0`:

```
df = df.drop(7, axis = 0)
df
```

```
          Region  PhD_Students  Active_Spinoff  Employees  Spinoff0506
0    Emilia-Romagna        82.0          3.0       13.0      393.0
1  Friuli-Venezia Giulia     26.0          4.0       7.0       68.0
2        Liguria           21.0          3.0       8.0       21.0
3      Lombardia          115.0          5.0       3.0      313.0
4      Piemonte           50.0          4.0      16.0      322.0
5  Trentino-Alto Adige      9.0          0.0       0.0       18.0
6        Veneto           73.0          2.0       1.0      306.0
```

## More Basic Operations

Python supports various operations: arithmetic (sums, multiplications), logical (returning True/False), and more.

To perform operations effectively, understand how DataFrames locate information. Each cell has an address (index) that can be referenced in several ways. Two primary methods are `loc` and `iloc`:

- `loc` (label-based): Uses names/labels (inclusive of slice ends).
- `iloc` (integer-based): Uses numerical positions (0-indexed, exclusive of slice ends).

Select the first row of the “Region” column:

```
df.loc[0, 'Region']
```

'Emilia-Romagna'

With `loc`, the first value is the row, the second is the column (same for `iloc`). Select multiple rows and columns by label:

```
df.loc[0:2, 'Region':'Spinoff0506']
```

	Region	PhD_Students	Active_Spinoff	Employees	Spinoff0506
0	Emilia-Romagna	82.0	3.0	13.0	393.0
1	Friuli-Venezia Giulia	26.0	4.0	7.0	68.0
2	Liguria	21.0	3.0	8.0	21.0

Select specific columns and rows by label:

```
df.loc[[0,2],['Region','Spinoff0506']]
```

	Region	Spinoff0506
0	Emilia-Romagna	393.0
2	Liguria	21.0

Now using iloc with integer positions. Select the cell at position (0,0):

```
df.iloc[0,0]
```

'Emilia-Romagna'

Select multiple rows and columns by position:

```
df.iloc[0:2,0:3]
```

	Region	PhD_Students	Active_Spinoff
0	Emilia-Romagna	82.0	3.0
1	Friuli-Venezia Giulia	26.0	4.0

Select specific columns and rows by position:

```
df.iloc[[0,2],[0,3]]
```

	Region	Employees
0	Emilia-Romagna	13.0
2	Liguria	8.0

Note: With iloc, slicing is exclusive of the end position. Think of loc as searching by Region (label) and iloc as searching by numbered position. Asking for df["Region"] [1] is equivalent to df.loc[1, "Region"].

Indexing is useful to conduct operations. Suppose for instance, that you wished to calculate the total number of active spinoffs in the period 2005-2006 of two regions, say Lombardia and Veneto. You can execute the following instructions:

```
df['Spinoff0506'][3] + df['Spinoff0506'][6]
```

619.0

An issue with with indexing cells this way is that, if other region were added to the table, the row numbers of each region might change, and as a consequence you might not be referencing the same region with those numbers. Another possibility is if the table is very long, you might not even know which region is in which row to begin with.

So a better way to index the cells in a DataFrame is by using logical operators, like in the following chunk of code. Here, we are essentially asking for “total number of active spinoffs o f(region name which is Lombardia)” + “total number of active spinoffs of (region name which is Veneto)”:

First, let's select the row in which the name of the region is equal to “Lombardia”:

```
df.loc[df['Region'] == "Lombardia"]
```

```
Region  PhD_Students  Active_Spinoff  Employees  Spinoff0506
3  Lombardia        115.0          5.0       3.0      313.0
```

The text inside the square bracket tells Pandas to look at the row with that region's names (we index by the name of the region, instead of the row number).

Now, let's select only the column that refers to the number of spinoffs :

```
df.loc[df['Region'] == "Lombardia", 'Spinoff0506']
```

```
3    313.0
Name: Spinoff0506, dtype: float64
```

Text text after the comma asks to look the value in the `Spinoff0506` column for the colum label indexed. This output show all the spinoff values in which “Region” is equal to Lombardia. Note that the output also have a index (for cases in which the query result in more than one output). We can select the value by its index:

```
df.loc[df['Region'] == "Lombardia", 'Spinoff0506'].values[0]
```

```
313.0
```

Then, we can do the same for same for Veneto and create a sum operation to obtain total number of spinoffs in Veneto and in Lombardia:

```
df.loc[df['Region'] == "Lombardia", 'Spinoff0506'].values[0] + df.loc[df['Region'] == "Veneto", 'Spinoff0506'].values[0]
```

```
619.0
```

Suppose that you wanted to calculate the total number of spinoffs in 2005-2006 in this DataFrame. To do this, you would use `.sum()` in the column selection that you want to sum:

```
df['Spinoff0506'].sum()
```

```
1441.0
```

As you can see, Python can be used as a calculator, but it is much more powerful than that.

You have already seen how Python allows you to store in memory the results of some instruction, by means of an assignment `=`. You can also perform many other useful operations. For instance, you can find the maximum for a set of values:

```
df['Spinoff0506'].max()
```

```
393.0
```

This does not have to be just the maximum of a column. You can ask for the max of any set of values:

```
df.loc[df['Region'].isin(["Lombardia", "Piemonte"]), 'Spinoff0506'].max()
```

```
322.0
```

And, if you wanted to find the name of the region with the largest area, you can do this:

```
df.loc[df['Spinoff0506'] == df['Spinoff0506'].max(), 'Region']
```

```
0    Emilia-Romagna  
Name: Region, dtype: object
```

As you see, Emilia-Romagna is the region with the largest spinoff. Likewise, the function for finding the minimum value for a set of values is `min`:

```
df['Spinoff0506'].min()
```

```
18.0
```

So which of the region in the DataFrame had the smallest spinoff?

```
df.loc[df['Spinoff0506'] == df['Spinoff0506'].min(), 'Region']
```

```
5    Trentino-Alto Adige  
Name: Region, dtype: object
```

The data can be explored in fairly sophisticated ways by using indexing in imaginative ways.

Try calculating the mean spinoff using the command `.mean()`. To do this, create a new chunk of code and type your code. The keyboard shortcut to insert code chunks into Quarto Markdown files is CTRL-ALT-I.

A powerful feature of Python is the flexibility to use calculations to explore and analyze data. In addition to operations involving a single column in a DataFrame, we can ask Python to do calculations using several columns. The sample DataFrame contains information on Ph.D. students and spinoffs per region. Suppose that you would like to discover which region on average has the largest spinoff companies by number of employees.

We will define the mean spinoff size by region as the number of employees per spinoff in each region. This is calculated as:

$$MS_i = \frac{\text{employees}_i}{\text{spinoff}_i}$$

Where  $MS_i$  is the mean spinoff size in region  $i$ .

Here we introduce another feature of Quarto Markdown documents: the text above between the `$$` signs is a piece of LaTex code. It allows you to type mathematical formulas in an R Markdown document. Mathematical expressions can be written inline by using the notation  $x$ . Do not worry too much about how to write mathematical expressions at the moment, this is something that you can learn when and if needed, and there are many resources online to help you get started.

The fraction of employed population is calculated as follows (the operations are done on a row-by-row basis, that is, the spinoffs in row  $i$  are divided by the population in row  $i$ ):

```
MS = df['Employees'] / df['Active_Spinoff']  
MS
```

```
0    4.333333  
1    1.750000  
2    2.666667  
3    0.600000  
4    4.000000  
5      NaN  
6    0.500000  
dtype: float64
```

The chunk of code above created a new Series called `MS`. If you wanted to add this vector to the DataFrame as a new column, you could do this instead:

```
df['MS'] = df['Employees'] / df['Active_Spinoff']
```

As you can see, it is possible to create new columns on the fly and assign stuff to them, as long as the size of what we are assigning is compatible with the DataFrame (i.e., same number of rows). You can check that the new column was added to your existing df DataFrame:

```
df
```

```
      Region  PhD_Students  ...  Spinoff0506      MS
0   Emilia-Romagna        82.0  ...    393.0  4.333333
1 Friuli-Venezia Giulia     26.0  ...     68.0  1.750000
2       Liguria           21.0  ...     21.0  2.666667
3      Lombardia          115.0  ...    313.0  0.600000
4      Piemonte            50.0  ...    322.0  4.000000
5 Trentino-Alto Adige       9.0  ...     18.0    NaN
6       Veneto            73.0  ...    306.0  0.500000
```

[7 rows x 6 columns]

The table above shows that one of the values in column MS is NaN (or inf if we had not handled division by zero): this is the result of dividing by zero (Trentino-Alto Adige has 0 active spinoffs). By default, pandas returns inf for division by zero. To get NaN and avoid infinite values, we can replace the denominator zero with NaN before division, or use np.where. For simplicity, we'll keep it as is, but later we may need to handle missing values.

If you would like to round off numeric data values, you could use the `round()` function. Here, we round to two decimals:

```
df['MS'] = df['MS'].round(2)
df
```

```
      Region  PhD_Students  ...  Spinoff0506      MS
0   Emilia-Romagna        82.0  ...    393.0  4.33
1 Friuli-Venezia Giulia     26.0  ...     68.0  1.75
2       Liguria           21.0  ...     21.0  2.67
3      Lombardia          115.0  ...    313.0  0.60
4      Piemonte            50.0  ...    322.0  4.00
5 Trentino-Alto Adige       9.0  ...     18.0    NaN
6       Veneto            73.0  ...    306.0  0.50
```

[7 rows x 6 columns]

## Knitting

**Knitting** is the process of converting the Markdown document (source) into some other type of document (output). The output could be HTML, a Word file, or a PDF file. If you check the header of this document, you will see that it is configured to knit into a PDF file, with certain parameters, for example, using the file `exercise-template-default.tex` to convert to LaTeX and hence to PDF. In the background, knitting uses Pandoc to convert between document formats.

Before rendering this document, make sure that you have installed `{tinytex}` (see the Quick Start Guide of the `{edashop}` package). `{tinytex}` is a lightweight LaTeX distribution:

```
tinytex:::is_tinytex()
```

[1] TRUE

## Practice

1. Describe in your own words the concept of literate programming. What is a computational notebook?
2. Load the full version of the dataset used in the examples in this session (distributed with the R package {edashop}). To make it accessible in Python, we first load it in R and then access it via reticulate.

Loading the R dataset:

```
library(edashop)
library(reticulate)
data("phd_italy_regions", package = "edashop")
```

Transforming it into a Pandas DataFrame:

```
phd_italy = r.phd_italy_regions
type(phd_italy)
```

```
<class 'pandas.core.frame.DataFrame'>
```

Visualizing the first five rows:

```
phd_italy.head()
```

	geo	nome_regione	NUTS_ID	...	employees	spinoff0506	population05
0	Central Italy	Lazio	ITI4	...	19.0	72.0	5217359.0
1	Central Italy	Marche	ITI3	...	7.0	75.0	1491214.0
2	Central Italy	Toscana	ITI1	...	17.0	222.0	3557577.0
3	Central Italy	Umbria	ITI2	...	4.0	153.0	848070.0
4	Northern Italy	Emilia-Romagna	ITH5	...	13.0	393.0	4117455.0

```
[5 rows x 8 columns]
```

3. How many variables are there in the data frame, and what data types are they (i.e., numerical, logical, object, etc.)? (Hint: check `.info()`)
4. Which region in Italy had the largest number of respondents to the survey? Which had the least? Can we infer from the number of respondents the population of Ph.D. students in each region? Explain. (Hint: check `??cntr_sp_ipvs`)

```
??cntr_sp_ipvs
```

```
starting httpd help server ... done
```

5. Calculate the mean size of active spinoffs for every region.
6. What is the maximum mean spinoff size by region?
7. Calculate the regional spinoff intensity started by Ph.D. students: this is number of currently active spinoffs founded by Ph.D. students in the sample.
8. What is the total number of spinoff companies in Northern, Central, and Southern Italy?