

# Getting Ready for your SW Carpentries Session



Felix Henninger, Sonja Grath, Reema Gupta,  
Malika Ihle and Luisa F. Jimenez-Soto



All material is Licensed under Creative Commons BY SA 4.0

July 23, 2024

# The Carpentries course

Dear Participant,

Welcome! We are happy you are interested in learning more about programming!

In general, Carpentries trainings, like this one, are designed to give you the basic tools and concepts for you to start to explore programming. You can use it for Data Analysis, programming, being lazy (let computers support your work!), and/or improving your work and research.

This course is called "Python flavored", because we will be teaching the basic parts: Unix and its shell, Git, and will focus on the programming language called python.

The structure of the course:

- Day 1: It contains two sections: Unix(Shell) and Git
- Day 2: Python: Programming and visualization

We will be doing mostly live coding sessions in order to give you an insight on how to work with them, but as any language: You need to practice it and if you do not understand, please ask. We will be at least 4 instructors and we are there to help you.

In order to make your experience easier, we have written this document for you to prepare and get the most of the training.

Please read it carefully, and if you have any questions feel free to ask us in the Help desk session scheduled for the day before the training.

Looking forward to meet you all,

your Instructors and helpers.

# Getting to know your computer

## Disclaimer

Your participation in the course requires that you download and install certain software and libraries maintained by an open community. If you decide to download and install them, we, the instructors, the OSC and the University (LMU) cannot be held liable for any damages that can occur in your computer. To our best knowledge and limitations, these programs and libraries have proven safe in previous trainings. This is the reason we recommend them and use them. However, each computer system is different. Make sure you take an informed decision before using them.

## To the point: Hardware and general software

It is very probable that you have worked with computers for at least a couple of years. Congratulations! However, for this course you will need to get to know your computer even better, specifically your keyboard.

**Before classes start**, we need you to **find out** the following tasks in your computer and keyboard

- **Your Terminal:** How to open your terminal (Mac OS and Linux OS) or PowerShell (Windows OS), (General information under <https://swcarpentry.github.io/shell-novice/>)

If you are a Windows user, please follow the instructions under [https://carpentries.github.io/workshop-template/install\\_instructions/#shell](https://carpentries.github.io/workshop-template/install_instructions/#shell) for Windows system,

MacOS and Linux are both Unix systems, so they have by default a shell. Just check the instructions in the same url mentioned above for your systems to verify you have what you need, and

- how to write the following characters:
  - Caret → `^`
  - Back slash → `\`
  - Forward slash → `/`
  - Square brackets → `[` and `]`
  - Curly brackets → `{` and `}`
  - Round brackets → `(` and `)`
  - Backtick → ```
  - Tilde → `~`
  - Pipe → `|`

- Pound or hash → #
- Ampersand → \$
- Underscore → \_
- Double quotes → " and "
- Single quote → ' and '
- Angle brackets → <and >
- Colon → :
- Semicolon → ;
- Comma → ,

All these characters will be used to communicate with your computer. The last thing you want is not being able to type the commands.

If you want more information about the types of keyboards used, here are some links and sources:

### Keyboard layouts (most common)

- American QUERTY: [https://en.wikipedia.org/wiki/QWERTY#/media/File:KB\\_United\\_States.svg](https://en.wikipedia.org/wiki/QWERTY#/media/File:KB_United_States.svg)
- QUERTZ German: <https://en.wikipedia.org/wiki/QWERTZ#/media/File:German-Keyboard-Layout.png>
- AZERT (France and Turkey): [https://en.wikipedia.org/wiki/AZERTY#/media/File:KB\\_France.svg](https://en.wikipedia.org/wiki/AZERTY#/media/File:KB_France.svg)

### Articles regarding Keyboards

- Contains the map of the distribution of keyboards in Europe and info about QUERTZ: <https://en.wikipedia.org/wiki/QWERTZ>
- American QUERTY: <https://en.wikipedia.org/wiki/QWERTY>
- French AZERTY: <https://en.wikipedia.org/wiki/AZERTY>

# Working with anaconda or miniconda

Conda is a package manager that allows us to use several versions of python in our computer without causing major conflicts.

In the instructions to the course, we gave you the option to install Anaconda or miniconda, two programs to use Conda. Here you have to choose: Anaconda is the easier way to interact with Conda, since it has the Anaconda Navigator, a Graphical User Interface (GUI) beginners friendly. However, *Anaconda is a very big program containing a lot of python libraries that might not be necessary for many people.* For Anaconda, follow [these instructions](#).

Your other option is to use Miniconda, the lighter version with the minimum content to start working with python. However, miniconda does not have a (GUI) like Anaconda does, and it requires installation of the individual packages using the command line interface. If you choose to install miniconda, you will find the information in this section. This includes what you need for installing the miniconda and how to install the libraries we will use for the training: Matplotlib, Pandas, and Jupyter Lab.

For the use of Jupyter Lab using the GUI, you can go to [Summary and Setup Lesson](#) from Software Carpentry.

## Installing miniconda

Depending on the operating system you have (MacOS, Windows, Linux), you need to download the program. Follow the instructions from the [official conda site](#) for the installation in your system.

For those of you who like it more visual, here are the some resources available in YouTube:

- For Windows 10: <https://youtu.be/oHHbsMfyNR4?si=WBb0zWVRgDbUynJh>
- For MacOS: <https://youtu.be/IRINDtw5xro?si=mfTCGE3-if7BCH0w>

## Using your installed miniconda to install libraries

Once you have installed your miniconda, you need to verify that the installation is working in your computer.

In your command line terminal / PowerShell type the following

---

```
1 $ conda --version
```

---

You should get something like `conda 23.1.0`. If you do not have the same version, type after the prompt (`$` or `>`) the following:

---

```
1 conda install conda
```

---

And accept the packages to be installed by typing `Y` or enter. You have now the most up-to-date version of conda.

For the course we need to use a couple of packages / libraries for data analysis. We are going to use the same method: interacting with conda using the command line, as we did for the update.

```
1 # For pandas and numpy
2 conda install anaconda::pandas
```

---

```
1 # for jupyter lab
2 conda install conda-forge::jupyterlab
```

---

```
1 # For matplotlib, for graphics
2 conda install conda-forge::matplotlib
```

Now you need to make sure you have everything you need for the course. For this you need to get the list of packages inside your conda environment called `base`. Just type after the prompt:

```
1 # Get the list of all packages inside your conda environment
2 conda list
```

Inside the list you should find the following (versions might differ a little bit):

- numpy (Version 1.26.4)
- pandas (Version 1.3.5)
- matplotlib (Version 3.9.0)
- jupyterlab (Version 4.2.3)

**Recommendation 1:** If you need help... we will be there for you!

If for some reason you get stuck or you do not feel confident enough about doing all this alone, do not despair! We have a Help desk (online) planned for **Wednesday, July 31st from 14:00 to 15:00** where we will try to answer any question you have about the installation process. Check your confirmation e-mail for the Zoom link.

# Cheatsheet Python

The following pages contain the different commands we will be using in class, on the second day. We want you to be able to keep this information for your future work, and maybe use it in class if you get lost (can happen!). This snippets are based on the course [Plotting and Programming with Python](#) created by the [Software Carpentry](#).

## Python Fundamentals

### Running and Quitting

---

```
1 print("hello world")
```

---

### Variables to store values

---

```
1 # Create the variable name and assign a value
2 age = 42
3 first_name = 'Ahmed'
4
5 # Use the print() function to display the values
6 print(first_name, 'is', age, 'years old')
7
8 # Variables can be use for mathematical operations
9 age = age + 3
10 print('Age in three years:', age)
```

---

### Access values by indexing

---

```
1 # Use the index to access one item in the variable
2 atom_name = 'helium'
3 print(atom_name[0])
4
5 # Slice to get a substring
6 atom_name = 'sodium'
7 print(atom_name[0:3])
```

---

## Data types and type conversion

---

```
1 # Declare some variables (assign variables to data)
2 first_value = 1.0
3 second_value = "1"
4 third_value = "1.1"
5
6 # Check the type of data stored under the variable name
7 print(type(first_value))
8 print(type(second_value))
9 print(type(third_value))
10
11 # Transform the data
12 print(int(second_value) + float(third_value))
```

---

## Built-in Functions and Help

---

```
1 # Obtain the maximum number in a series
2 max(1, 2, 3)
3
4 # Obtain the minimum
5 min(1, 2, 3)
6
7 # evaluate the length of the string "hello"
8 len("hello")
9
10 # Check the information about the function len()
11 help(len)
```

---

## Lists

---

```
1 # Create a list
2 odds = [1, 3, 5, 7]
3
4 # print elements of the list using indexes
5 print('first element:', odds[0])
6 print('last element:', odds[-1])
7
8 # modify the list
9 odds.append(9)
10 print('odds after append:', odds)
11 odds.extend([11, 13, 15])
12 print('odds after extend:', odds)
13 del odds[0]
14 print('odds after delete:', odds)
```

---



## Loops

---

```
1 # create the list
2 squares = [1, 4, 9, 16]
3
4 #loop over each item of the list using index for selection
5 for i in range(len(squares)):
6     squares[i] = squares[i] + 1
7     print(squares)
```

---

## Conditionals

---

```
1 num = 37
2
3 # check if the value stored under num is greater than 100
4 if num > 100:
5     print("greater")
6 else:
7     print("not greater")
```

---

## Writing Functions

---

```
1 # declare and define the function to transform fahrenheit to celcius
2 def fahr_to_celsius(temp):
3     return ((temp - 32) * (5/9))
4
5 # apply the function and print the value
6 print(fahr_to_celsius(32))
7 print(fahr_to_celsius(212))
```

---

## Libraries

---

```
1 # import the library called math
2 import math
3
4 # use the function from math to access the constant pi and print it
5 print("pi is", math.pi)
6
7 # Use the function cos() from math library to calculate the cosine of pi, and
8   print the answer
9 print("cos(pi) is", math.cos(math.pi))
```

---

# Data Analysis using Pandas library

## Reading Tabular Data into DataFrames

---

```
1 import pandas as pd
2
3 # use the fuction read_csv() from the library pandas (pd) to read the csv file
  with the data
4 data = pd.read_csv('gapminder_gdp_oceania.csv')
5
6 #print the dataframe resulting from using the pandas function read_csv()
7 print(data)
```

---

## Pandas DataFrames

---

```
1 # print on the console the columns inside the dataframe under the variable data
2 print(data.columns)
3
4 # get an idea of the information stored in the DataFrame and its organization
5 print(data.describe())
```

---

## Looping Over Data Sets

---

```
1 import pandas as pd
2 #Create a list containing the names of files you want to open and read in a
  dataframe
3 filenames = ['data/gapminder_gdp_africa.csv', 'data/gapminder_gdp_americas.csv']
4
5 #Create an empty list where all dataframes will be placed
6 dataframes = []
7
8 # Loop over the 'filenames' and read each in a dataframe
9 for f in filenames:
10     dataframes.append(pd.read_csv(f))
11     print(f, pd.read_csv(f).describe())
```

---

## Plotting and Matplotlib library

---

```
1 # import the package for plotting (matplotlib.pyplot)
2 import matplotlib.pyplot as plt
3
4 # create a barplot using the values inside the dataframe under "data" variable
5 data.plot(kind='bar')
6
7 # show the plot
8 plt.show()
```

---

# Advanced extras

## Errors and Exceptions

---

```
1 def favorite_ice_cream():
2     print('My favorite ice cream is', ice_cream)
3
4 favorite_ice_cream()
```

---

## Defensive Programming

---

```
1 def normalize_rectangle(rect):
2     if rect[0] * rect[1] < 0:
3         raise ValueError('Invalid rectangle')
4     if rect[0] < 0:
5         rect[0] = -rect[0]
6     if rect[1] < 0:
7         rect[1] = -rect[1]
8     return rect
```

---

## Variable Scope

---

```
1 pressure = 103.9
2
3 def adjust(t):
4     pressure = 0.9 * t
5     return pressure
6
7 print(adjust(pressure))
8
9 print(pressure)
```

---

## Programming Style

---

```
1 # Good variable names
2 count = 10
3 name = "temp_data"
4
5 # Avoid magic numbers
6 MAX_TEMP = 100
7 if temperature > MAX_TEMP:
8     print("too hot")
```

---

# Open Science and FAIR principles

Having a computationally reproducible workflow is part of making science more efficient, transparent, and reliable. Other research practices can enhance the credibility of science, and all of them are often packaged into what is called ‘open research practices’. With the [Carpentries](#), the tools you learn here will help you to follow the principles of Open science and FAIR data.

## Open Science

How many of you have tried to reproduce an experiment you have read and it does not seem to work? Some of us have. And many times. These “reproduction of experiments” activity has wasted weeks of work and materials just to find out that the results and experiments described together with the protocols published do not bring the same results.

In 2005, a publication by John P.A. Ioannidis [1], discussed openly some of the first problems regarding reproducibility problems. This time it was for Behavioral Research (Psychology and co.), however, the same problems were later made visible in many other areas of science. A more mature view of the problem has led to defining the reasons, “symptoms” (including P-Hacking) and possible “cure” for the reproducibility crisis [2]. The Open Science movement started.

If you want to read more about what Open Science means, there are excellent resources in the [LMU Center for Open Science](#), teaching materials in the [Open Science book](#), or very nice people willing to help in the [OSIM \(Open Science Initiative in Medicine\)](#)

## FAIR principles

FAIR Part of the Open Science ideals is to be able to share the data. The standards for data and metadata (data about how data was obtained) under Open Science principles are defined under the FAIR acronym:

- F** Findable: The data has a unique identifier and the data complementing it (metadata) is descriptive enough that the identifier can be associated with the studies.
- A** Accessible: We can access the data with standard communication protocol.
- I** Interoperable: The data is in a standard format that can be read by all without the need to buy a software for it, and metadata have such standards (called ontologies) that the dataset can be merged with another.
- R** Reusable: The metadata contains all details about the data and its origin, so that it can be used by other people. The license has to allow this as well.

If you want more detailed information, and how to implement them go to the publication from Mark D. Wilkinson et al, “The FAIR Guiding Principles for scientific data management and stewardship” [3]

## Sources and Methods used for this guide

The code you find in this guide was extracted from the original material from the SW Carpentries training using ChatGPT 4o, and modified for understanding and simplicity.

The logos are from Open Sources available in Wikipedia under Creative commons or from the organizations themselves.

If you find an error, or a link is not working, please let us know so we can improve the manual. We are happy to get constructive feedback.

# Bibliography

- [1] J. P. Ioannidis, “Why Most Published Research Findings Are False,” *PLOS Medicine*, vol. 2, p. e124, jul 2005.
- [2] M. R. Munafò, B. A. Nosek, D. V. Bishop, K. S. Button, C. D. Chambers, N. Percie Du Sert, U. Simonsohn, E. J. Wagenmakers, J. J. Ware, and J. P. Ioannidis, “A manifesto for reproducible science,” *Nature Human Behaviour* 2017 1:1, vol. 1, pp. 1–9, jan 2017.
- [3] M. D. Wilkinson, M. Dumontier, I. J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J. W. Boiten, L. B. da Silva Santos, P. E. Bourne, J. Bouwman, A. J. Brookes, T. Clark, M. Crosas, I. Dillo, O. Dumon, S. Edmunds, C. T. Evelo, R. Finkers, A. Gonzalez-Beltran, A. J. Gray, P. Groth, C. Goble, J. S. Grethe, J. Heringa, P. A. t Hoen, R. Hooft, T. Kuhn, R. Kok, J. Kok, S. J. Lusher, M. E. Martone, A. Mons, A. L. Packer, B. Persson, P. Rocca-Serra, M. Roos, R. van Schaik, S. A. Sansone, E. Schultes, T. Sengstag, T. Slater, G. Strawn, M. A. Swertz, M. Thompson, J. Van Der Lei, E. Van Mulligen, J. Velterop, A. Waagmeester, P. Wittenburg, K. Wolstencroft, J. Zhao, and B. Mons, “The FAIR Guiding Principles for scientific data management and stewardship,” *Scientific Data* 2016 3:1, vol. 3, pp. 1–9, mar 2016.