

Introduction to Data Analysis with Python

Cecilia Graiff

February 5, 2026

ALMAnaCH, Inria Paris

cecilia.graiff@sciencespo.fr

Agenda

Introduction to coding workflows

Collaborative coding and version control

Virtual Environments

How does Python actually work?

Git and GitHub Basics

Sample Workflow

Understanding your computer

Introduction to coding workflows

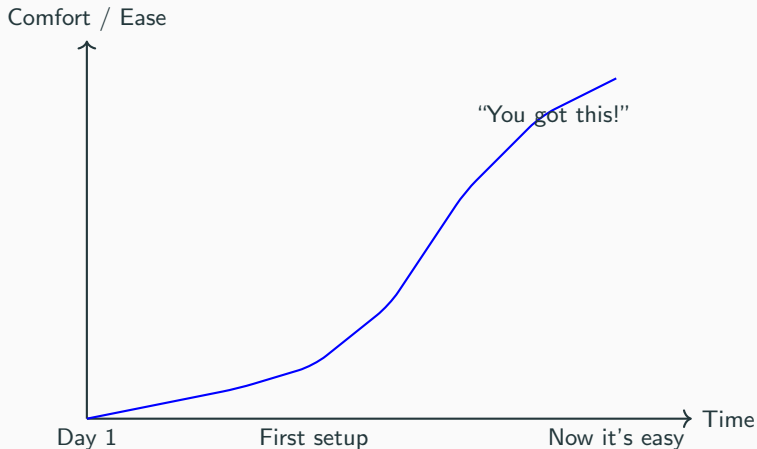
- The easiest option is **Colab Notebooks**
 - You only need a Google account (the SciencesPo one works!)
 - You do not need to install anything
 - You can collaborate dynamically, like on GoogleDocs

On Colab Notebooks, you also have **access to GPU**, which facilitates heavier computations (not necessary for this course).

Real-world pipeline

- Colab Notebooks not suitable for more complex projects
- Therefore: learning a **data scientist workflow**
- "Entrance fee": the setup can be confusing!

Your Learning Curve



Collaborative coding and version control

Coding and version control

- You write code **locally** on your computer.
- You use Visual Studio Code, which is **an editor**: think of it like a text editor, an app that you can use to **open, edit, and share your files**
- You track your work with **Git**.
- You share and collaborate through **GitHub**.

Coding and version control

- You write code **locally** on your computer.
- You use Visual Studio Code, which is **an editor**: think of it like a text editor, an app that you can use to **open, edit, and share your files**
- You track your work with **Git**.
- You share and collaborate through **GitHub**.

To sum up

Think of this as writing a report with your own research setup (Python), saving versions (Git), and uploading to a shared folder (GitHub).

What is Python?

- Python is a **high-level programming language**.
- It is easy to read and write, mimicking English grammar.

What is Python?

- Python is a **high-level programming language**.
- It is easy to read and write, mimicking English grammar.
- Widely used for:
 - Data analysis and statistics
 - Web development
 - Automation and scripting
 - Machine learning and computational social science
- Python uses **modules and libraries** to extend functionality, e.g., `numpy`, `pandas`, `matplotlib`.

Python Code: Example

```
# Print hello world
print("Hello, Python!")

# Simple computation
x = 5
y = 10
print("Sum:", x + y)
```

How a Python Project Works

- A Python project is usually a folder containing code and data.

How a Python Project Works

- A Python project is usually a folder containing code and data.
- Typical structure:
 - `src/` – repository containing your code; for example, it could contain `preprocessing.py`, `visualization.py`, and `linear_regression.py`
 - `data/` – folder for datasets
 - `requirements.txt` – lists external libraries your project uses, to ensure reproducibility
 - `README.md` – project description and instructions

How a Python Project Works

- A Python project is usually a folder containing code and data.
- Typical structure:
 - `src/` – repository containing your code; for example, it could contain `preprocessing.py`, `visualization.py`, and `linear_regression.py`
 - `data/` – folder for datasets
 - `requirements.txt` – lists external libraries your project uses, to ensure reproducibility
 - `README.md` – project description and instructions
- Use **virtual environments** to manage dependencies:
 - Keeps project libraries separate from system Python
 - Allows reproducing your setup

Virtual Environments

What is a Virtual Environment?

- A **virtual environment** is like a private workspace for your project.
- It is a **folder** where you tell Python to download specific packages
- Often, two projects require different packages, or different versions of the same package: hence the **need for a virtual environment**
- It contains its own copy of Python and packages.
- This prevents version conflicts between projects.
- **Reproducibility is fundamental**: to execute your code, people need to **install the same versions of the same packages**

Virtual Environments

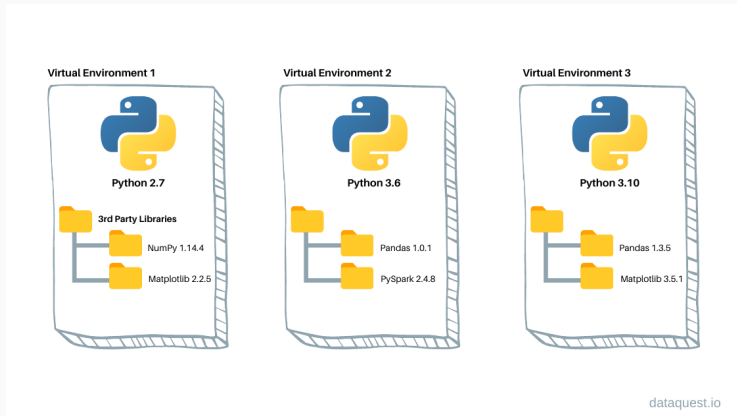


Figure 1: Virtual Environments in Python

Source: [Dataquest](https://dataquest.io)

Manage a Virtual Environment

1. Activate the virtual environment

```
cd /path/to/project
```

```
python -m venv venv
```

Manage a Virtual Environment

1. Activate the virtual environment

```
cd /path/to/project
```

```
python -m venv venv
```

2. Activate the virtual environment

```
# Windows
```

```
venv\Scripts\activate
```

```
# macOS/Linux
```

```
source venv/bin/activate
```

Manage a Virtual Environment

1. Activate the virtual environment

```
cd /path/to/project  
python -m venv venv
```

2. Activate the virtual environment

```
# Windows  
venv\Scripts\activate
```

```
# macOS/Linux  
source venv/bin/activate
```

Ideally, you **create** a virtual environment for a project only once, and you **activate** the virtual environment any time you work on the project.

Install Packages

- Python relies on **libraries** (such as pandas, matplotlib or numpy) and **packages** (such as matplotlib.pyplot)
- Those packages need to be **installed** in the virtual environment (**only once**):

```
pip install pandas
```

```
pip install numpy
```

Install Packages

- Python relies on **libraries** (such as pandas, matplotlib or numpy) and **packages** (such as matplotlib.pyplot)
- Those packages need to be **installed** in the virtual environment (**only once**):

```
pip install pandas
```

```
pip install numpy
```

- Any time you run a **notebook**, you need to **import** the libraries that you will use:

```
import pandas as pd
```

```
import numpy as np
```

Install Packages

- Python relies on **libraries** (such as pandas, matplotlib or numpy) and **packages** (such as matplotlib.pyplot)
- Those packages need to be **installed** in the virtual environment (**only once**):

```
pip install pandas  
pip install numpy
```

- Any time you run a **notebook**, you need to **import** the libraries that you will use:

```
import pandas as pd  
import numpy as np
```

Note the important difference: the first command is for your **computer** (you are downloading something), the second one is **following Python syntax** and refers to your Python script.

How does Python actually work?

Python interpreter

- When you install Python, an **interpreter** is downloaded: this interpreter **understands what you write in Python** and **runs it** (= makes the execution possible)
- Different **operating systems (OS)** call this interpreter differently
- Often, you need to add Python to the **Path variables** if you want to be able to use the standard commands

- Possible commands:
 - `python` usually refers to the latest Python, but if you have both `python3` and `python2`, it can refer to Python 2
 - `python3`: usually safest option
 - `py`: on Windows, Python comes with an interpreter called `py`. Using this commands solves eventual problems with the `PATH` variable (if you code on a regular basis, I recommend adding Python to `PATH`, but otherwise, it can be fine)

- Possible commands:
 - `python` usually refers to the latest Python, but if you have both `python3` and `python2`, it can refer to Python 2
 - `python3`: usually safest option
 - `py`: on Windows, Python comes with an interpreter called `py`. Using this commands solves eventual problems with the `PATH` variable (if you code on a regular basis, I recommend adding Python to `PATH`, but otherwise, it can be fine)

Please try to understand what your setup is, and what you need to type. You might have to change my code examples, where I usually use `python`.

Git and GitHub Basics

Why Use Git and GitHub?

- **Git** tracks versions of your code (like “Save As” with history).
- **GitHub** hosts those versions online.
- Together, they make your projects **reproducible and shareable**.

Git vs GitHub

- Git is a **tool for version control**
 - Therefore: you need to install Git (did at the beginning of the course)
- GitHub is a **website built around Git**. It is not the only existent one, but the most widely used.
 - Therefore: No need to install anything (although a tool called GitHub Desktop exists: it is a downloadable interface, but not required for this class: due to its limitation, most people use only GitHub).

Git: Reproducibility and Version Control

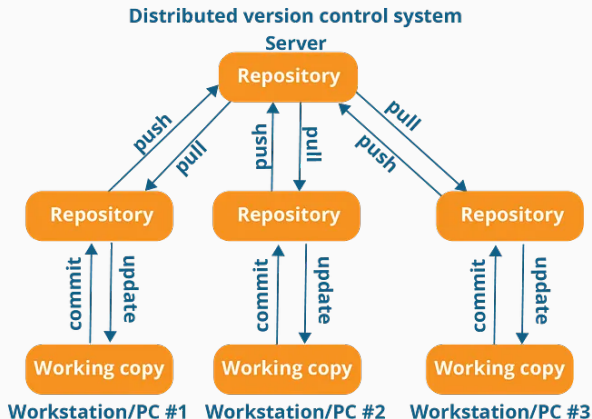


Figure 2: A simplified diagram of version control in Git.

Basic Git Workflow

1. Clone the repository:

```
git clone https://github.com/username/project.git
```

2. Track and save changes:

```
git add .  
git commit -m "Describe your change"
```

3. Upload (push) to GitHub:

```
git push
```

Basic Git Workflow

1. Clone the repository:

```
git clone https://github.com/username/project.git
```

2. Track and save changes:

```
git add .  
git commit -m "Describe your change"
```

3. Upload (push) to GitHub:

```
git push
```

Summary: Activate your environment, work locally, then commit and push.

Sample Workflow

Example Workflow

1. Open terminal and navigate to your project folder.
2. Activate your virtual environment.
3. Run or edit your notebooks/scripts.
4. Save and commit changes with Git.
5. Push to GitHub to back up and share.

Example Session:

```
cd myproject
```

(or you open the project directly in VSC and use the VSC terminal)

Example Session:

```
cd myproject
```

(or you open the project directly in VSC and use the VSC terminal)

```
source env/bin/activate
```

(install packages, modify your files, write code, etc)

```
git add .
```

```
git commit -m "Updated analysis"
```

```
git push
```

Understanding your computer

Terminal vs Python commands

- Terminal: you write here what you need to tell **to your computer**
 - Example: git commands, pip install
 - Those are not commands specific to the Python file/Jupyter notebook you are using: When you install a package, you install it **in your virtual environment**, that you can use for any Jupyter notebook you want to create
 - The terminal can be opened with the dedicated app, or in VSC (Click on "Open terminal" in the menu, and it will be displayed in the lowest half of the screen)
 - Important: for pip install, you can have a **shortcut** by talking to your computer **directly from a notebook**, if you **add an exclamation mark before the command**: `!pip install pandas`

Terminal vs Python commands

- Python: Here you write commands specific to your Python file/Jupyter notebook, **with the Python syntax**
 - Example: import packages, functions for data analysis

Terminal vs Python commands

- Python: Here you write commands specific to your Python file/Jupyter notebook, **with the Python syntax**
 - Example: import packages, functions for data analysis

Python is a programming language and as such has a **syntax**, which is different from the one of terminal commands.

Questions?