**Máster en Big Data. Tecnología y Analítica Avanzada (MBD).**

**Introducción al Análisis Estadístico con Lenguajes de Programación para Machine Learning (IAELPML). 2023-2024.**
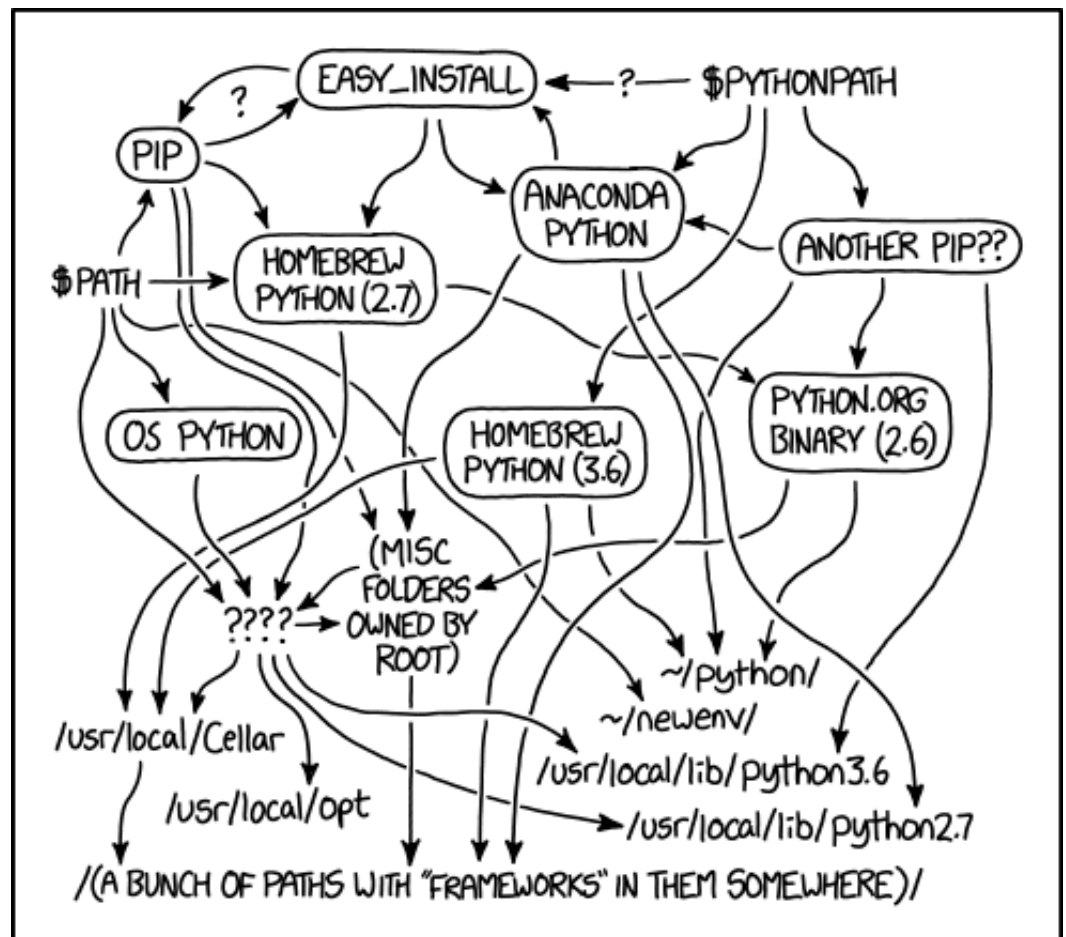
**Software Setup, MacOS Version**

## Before you proceed

- Please, read this document completely before starting your setup, so as to get a general idea of what is coming. In particular, please read the final *Alternative Setups* section. If you have any doubts, we are happy to help.
- **Be careful!** Another reason to read this document in advance is that we have seen too often users ruining their Python setups because of a rushed decision. Wrecking your system can be a matter of minutes, getting it back can take you hours or even days. Don't let the cartoon below be you!! When in doubt, please talk to us first.



*Image from XKCD (https://xkcd.com/1987/)*

# Folder structure for the course

- You should create a **specific folder** for this course under your user folder. Under MacOS a typical location could be something like (the ellipsis represents possible intermediate folders):

  ```
  /Users/your_user_folder/.../MasterBD/IAE/
  ```

- **The folder that you have created will always be referred to in these notes as the Course Folder**. In the above example we named it `IAE` but you can give it any name you want. As usual, it is a good idea to keep it short and meaningful (and free of spaces, unusual characters, etc.)
- We will be using Git and [GitHub (https://github.com)](https://github.com) repositories for the Course Notes and Assignments. In particular, the following repository will be used to share all the common course material:
  [https://github.com/mbdfmad/IAELPML (https://github.com/mbdfmad/IAELPML)](https://github.com/mbdfmad/IAELPML)
  In the next section we will tell you how to get a local copy of this repository to your computer. For the time being you can browse the contents of the repository, but do not download it yet! This repository **will always be referred to in these notes as the Common Repository**. Soon you will be also using your *personal repository* and it is important to get the distinction and name conventions right from the beginning.

# Software Setup and First Steps

## Installing Homebrew

- [Homebrew (https://brew.sh)](https://brew.sh) is a *package manager* for macOS (or Linux), widely used to install software which is harder to install or unavalaible in other channels, such as the Apple App Store for MacOS. Many programming languages and utilities are easy to install using Homebrew.
- To install Homebrew itself you just need to foloow the simple instructions in the *Install Homebrew* right at the top of the HomeBrew web page. But **before doing that** we highly recommend reading through the following section, called *What Does Homebrew Do?*. And in case you have any doubts please ask us before proceeding.

## Installing Git

**Important Note:** If you already have a working and recent setup of Git (current version is 2.37), skip to the [*Installing Homebrew Python*](#) subsection. If you do not know, open a terminal and try to run this command:

```
git --version
```

If you receive an error message then run:

```
brew install git
```

This will take some time. After the installation, try running the `git --version` command again.
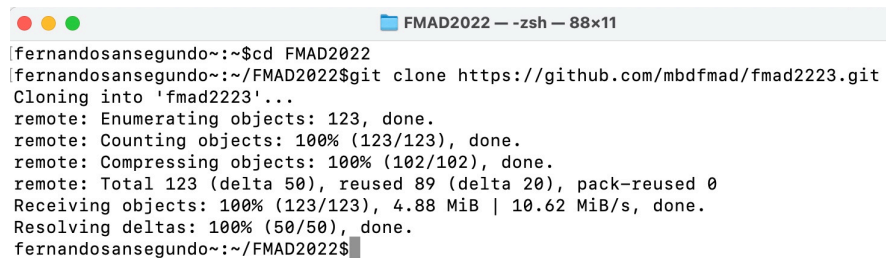
## Create a GitHub account

- **Important Note:** even if you already have a GitHub account you need to have an account that is linked with your university email (the one ending in `@alu.comillas.edu`). You can create that account using [this link (https://github.com)](https://github.com). Make sure to log out first if you were already logged with a *non-university-associated* account. And we recommend [reading this advice first (https://happygitwithr.com/github-acct.html#username-advice)](https://happygitwithr.com/github-acct.html#username-advice).

## Clone the *common repository*

- Using the Terminal navigate to the *Course Folder*. You can get there using the `cd` command. **Make sure that you have changed into the right folder before going further!!** You can use the `pwd` command to check your current location. Ask for help if you run into trouble.
- In that folder run the following command:

  ```
  git clone https://github.com/mbdfmad/IAELPML.git
  ```

  Git will begin cloning (downloading) the contents of the repo to a folder named `IAELPML` in your computer. The process looks like this (the prompt and number of cloned objects will surely be greater in your case):

  

  The `IAELPML` folder will then be a subfolder of the *Course Folder*. You should consider this `IAELPML` folder as a *read-only* folder. You will soon see how to update its contents as the course progresses, but **you are not expected to directly modify the contents of the `IAELPML` folder** (and you shouldn't). You can *copy* any object from this folder (e.g. the code examples that we will provide) to other locations in your computer if you want to make changes. Finally, if you are a seasoned Git user you can opt for a fork of the repo instead of cloning it.

## Keeping your local copy of the *common repository* updated.

- If you are used to auto syncing tools such as Dropbox, OneDrive, ICloud Drive and the like, you must keep in mind that **Git repos are not auto-syncing**. To keep the information in your repos updated you have to manually update them (both ways, from your computer to GitHub and viceversa).
- **This will be very important for our workflow in the lectures of the course.** So make sure to get it right, and ask for help soon if you need it. Every work session for this course will begin with the following ritual:
  - Open a *Terminal*.
  - Navigate to your local `IAELPML` folder. Do not confuse this with its parent folder, the *Course Folder*; in my case the parent course folder is called `IAE`, but in your case the name could be different. The repository folder name should however be always `IAELPML`.
  - Run `git status` to check that everything is ok.
  - Run `git pull origin main`. Depending on the update status of your local copy of the repo you will either download the new contents from the repo at GitHub or you will receive a message saying that you are `Already up to date`.
  - If you are going to modify any of the files in `IAELPML` and you want to keep the changes, **make sure to copy the files first to a different folder outside of the common repository.** Just create as many additional subfolders of the course folder as you see fit, but make so outside of the common repository.

  Don't worry too much, however. You cannot break anything, since you won't have write permissions for this GitHub repo. Worst case scenario, you can (backup your local files and) delete the `IAELPML` folder and clone it again to get the latest contents.

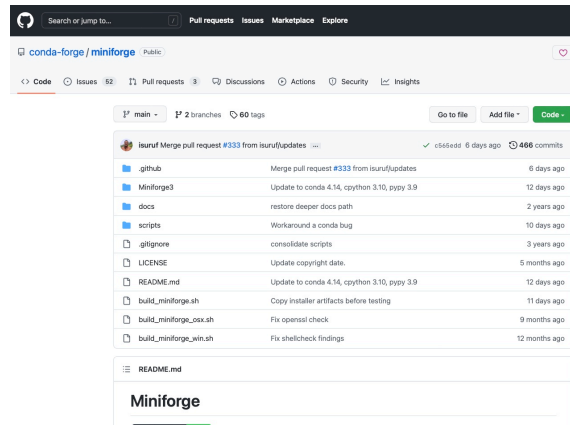## Setting up a local Git configuration for the Common Repository

- Using still the terminal, change into the `IAELPML` folder and set your local Git configuration with these commands (the first one is the folder change from the *Course Folder* to `IAELPML`):

```
cd IAELPML
git config --local user.name
'replace_this_with_your_GitHub_username'
git config --local user.email
'replace_this_with_your_university_email'
git config --global init.defaultBranch main
```

- By the way, if you are doing this setup for the first time, you can now close the *Git Bash* terminal window.

## Installing Python with HomeBrew

**Important Note:** If you already have a working setup of Python using conda for package and environment management (with Python 3.8 or later) then skip to the *Create a Conda Environment for the Course* subsection. If your Python version is older: get updated, people!

- Begin by navigating to the Miniforge GitHub repository (https://github.com/conda-forge/miniforge):



I suggest you at least take a look at the first sections. After that, if you skip the Downloads and scroll down to the `Install` section, there is a subsection called `Homebrew` (https://github.com/conda-forge/miniforge#homebrew). There you will see that all you need to do is open a terminal and run:

```
brew install miniforge
```

Wait for `HomeBrew` to finish setting the miniforge flavor of Python for you. Then in the same terminal execute

```
python
```

and you should see output like this:

```
Python 3.10.5 | packaged by conda-forge | (main, Jun 14 202
2, 07:07:06) [Clang 13.0.1 ] on darwin
Type "help", "copyright", "credits" or "license" for more i
nformation
>>>
```

The `>>>` line is the Python prompt waiting for you. But we will not be using Python this way by now. Therefore simply type `exit()` in the Python prompt and execute it to return to the regular prompt of the MacOS terminal.

# Create a Conda Environment for the Course

- A conda environment allows us to organize our work in projects, keeping each project dependencies isolated from the rest of your work. It is also the best way to avoid running into the potentially catastrophic situation in where you have several Python versions and they conflict to [create a mess (https://xkcd.com/1987/)](https://xkcd.com/1987/). We urge you to **always use environments** to organize your work in Python. You can learn more about the use of Conda environments [in this link (https://docs.conda.io/projects/conda/en/latest/user-guide/tasks/manage-environments.html#)](https://docs.conda.io/projects/conda/en/latest/user-guide/tasks/manage-environments.html#). And to learn more about environments in Python in general, you can check these two links:
    - [A Guide to Python's Virtual Environments (https://towardsdatascience.com/virtual-environments-104c62d48c54)](https://towardsdatascience.com/virtual-environments-104c62d48c54)
    - [The Definitive Guide to Conda Environments (https://towardsdatascience.com/a-guide-to-conda-environments-bc6180fc533)](https://towardsdatascience.com/a-guide-to-conda-environments-bc6180fc533)
- We will be using a Conda environment named `IAELPML` for this course. To create the environment we will keep working in the terminal (make sure that you exited Python in the previous step, you must be in a regular terminal prompt for this step to work).
- In that prompt execute the command

    ```
    conda create --name IAELPML
    ```

    and after some moments something like this should happen (with IAELPML instead of fmad in the following screenshots):

    ```
    (base) fernando@192 ~ % conda create --name fmad
    Collecting package metadata (current_repodata.json): done
    Solving environment: done

    ## Package Plan ##

      environment location: /Users/fernando/miniforge3/envs/fmad


    Proceed ([y]/n)? y

    Preparing transaction: done
    Verifying transaction: done
    Executing transaction: done
    #
    # To activate this environment, use
    #
    #     $ conda activate fmad
    #
    # To deactivate an active environment, use
    #
    #     $ conda deactivate

    Retrieving notices: ...working... done
    (fmad) fernando@192 ~ %
    ```

    For the time being ignore the warning about updating Conda if it appears (we will deal with it in the next step) and just answer *y* and hit *Enter*.

## Update conda

- To update conda (if the warning appeared in the previous step) run this command:

```
conda update –n base –c defaults conda
```

As the following picture illustrates conda will present you with a summary of the update operation. The picture is from Windows but the MacOS version is almost identical. Answer *y* and hit *Enter*.



Then conda will download the required software to autoupdate. The expected result looks like this:

## Activate the environment

- When we created the  IAELPML  environment in a previous step, Conda showed us the activation command:

  `conda activate IAELPML`

  Run this command now. This will activate the newly created environment `IAELPML`. Note that the name of the active environment appears in parenthesis at the beginning of the prompt line (again, here you will see `IAELPML` instead of `fmad`):

  <img src="fig/004-AnacondaEnvironmentCreate03.png" width=30%/>

  **Make sure that the environment is correctly activated before proceeding!!** If you have any problem ask for help.

## Install the first library

- The first library that we will install is [NumPy (https://numpy.org)](https://numpy.org), which provides the foundation for much of the work we will doing in this course. We will discuss NumPy in coming sessions of the course. But to use it we first need to install it in our `IAELPML` environment. Execute this command now:

  ```
  conda install --name IAELPML numpy
  ```

  A summary of the libraries that need to be downloaded for the NumPy setup will appear

  

  followed by a rather lengthy list of the libraries about to be installed (in most cases these two lists are almost equal)

  

  Answer yes to make the setup begin and wait patiently while NumPy and all of its requirements are downloaded and installed. The proccess should end like this:

  

## Install aditional libraries

- NumPy is just the first of the libraries that we will use in the course, but many others will also come into play. To name a few of the most relevant:
    - [NumPy (https://numpy.org)](https://numpy.org)
    - [Pandas (https://pandas.pydata.org)](https://pandas.pydata.org)
    - [MatPlotLib (https://matplotlib.org)](https://matplotlib.org)
    - [Scikit-Learn (https://scikit-learn.org/stable/)](https://scikit-learn.org/stable/)

    Take a few moments to visit the homepage for each of these libraries.

**Exercise:** install all these libraries to the `IAELPML` environment. To install several libraries you don't need to repeat `conda install` for each library; you can use a single install command and separate the names of the libraries with a space.

# Install the Jupyter Notebook & Jupyter Lab

- In this course we will talk to Python using the onterface provided by *Jupyter Notebooks*. The [Jupyter Project (https://jupyter.org)](https://jupyter.org) provides a general purpose development environment which is well suited for Python (but not limited to it). It will play a central role in our work with Python, but it installs (to the IAELPML environment) just like any other library.

  ```
  conda install --name IAELPML jupyter
  ```

- When the setup ends execute the following two commands in order. The first one changes your current directory to your user folder in Windows. The second command opens the Jupyter Notebook in your *default Web Browser*:

  ```
  cd %HOMEPATH%
  jupyter notebook &
  ```
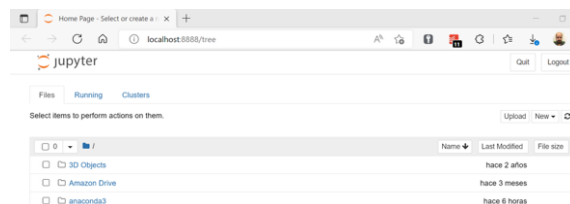
  Like this:



  *Note:* We recommend to always start Jupyter from a folder high enough in your folder structure: you won't be able to easily move up from the starting point later. Make sure, however, that you are at the very least, at the *Course folder* level.

  In a few moments your default *Web Browser* will open to a new page that shows the entry point for the running Jupyter server that we have just started:



You can use that page to navigate your folder structure until you get to the *Course Folder*. Once you are seeing the contents of that folder, open the 01a_test_your_setup.ipynb file and follow the instructions there. If errors appear try to go over the setup procedure again and if you are still having trouble ask for our help.

### Final Remarks

- The terminal window that you used to start Jupyter should remain open during all the work session. **Be careful:** Closing it can make the server crash and you may lose part of your work.
- Remember to make external copies of the files in `IAELPML` before modifying them.
- When you are done with the setup checks use the `File` menu and then `Close and Halt` the notebook. If your browser asks for confirmation that you really want to close the page: don't worry, you can safely close it.
- Every time you finish working with Jupyter you should go back to the `Home Page` of the Jupyter server (the one which shows your folder structure) and use the `Quit` button on the top right side of that window to shutdown Jupyter. A popup window will appear to inform you that you should close that browser tab.
- Return to the terminal that you used to initiate Jupyter. You should now see some messsage confirming that the server is shut down. If you are done working in this Python environment, issue the command:

  ` conda deactivate`

  The prompt will change to indicate that you are back to the `base` conda environment. You can now close the terminal window as well if you wish.

## Jupyter extensions

- If we need to install the extensions, please run this command:

  `conda install -n IAELPML jupyter_contrib_nbextensions`