

Practical work 01 – 19th of septembre 2017

Let's get started

The main objective of this Practical Work (PW) is to set all things up regarding your computer environment. We then have some exercises to wrap things up regarding what was said in the chapter 1 on **Fundamentals on Machine Learning**.

In this class we are going to use Python as programming language to perform some of the PW. So, this week's PW will include setting up Python on your computer and getting familiar with this language.

Python was released in 1991 and initially developed by Guido van Rossum (Holland). Python got inspiration from Modula 3, is object-oriented and can be easily extended. Good libraries are available to work with data, to visualise it and to model it with machine learning. This is why we have chosen to use Python in this class.

You have basically 3 ways to work with Python :

- a) Using a Python console in interactive mode. When Python is installed on your computer, just open a terminal and launch python.
- b) Using a text editor or an IDE to develop scripts of functions (i.e. saving ".py" files), as you would do it in Java or C. A good IDE is Pycharm from Jet Brains.
- c) Using iPython notebooks via a web browser. The concept is similar to Mathematica notebooks where you can mix code and text inputs (using Markdown) in so-called cells.

For this class we recommend to work with iPython Notebooks (option c) which is probably the best way to give us back your practical works. But you can also give back pure Python files together with a little report if you prefer (option b).

Warning : Which version of Python to use? There are now two "concurrent" versions of Python, versions 2 and 3. The Python world is now migrating strongly towards using Python 3. Unfortunately, Python 2 and Python 3 have some minor incompatible differences in language syntax and some modules and libraries did not move yet to 3. The choice is up to you. When we provide solutions to the exercises, we will let you know which version of Python has been used.

For this PW, **there is no need to return any report**. From next week, we will assume that your environment is all set up and that you have gotten familiar with Python.

Exercise 1 Moodle subscription

Register on moodle at <https://moodle.msengineering.ch>. To find the class navigate to Home → Lausanne → Approfondissement technico-scientifique (Type T) → S1 - 2017-2018 → T-MachLe - Machine Learning. Use subscription key : `tmachle1718`.

Most of the time, you'll need to submit your PW reports for the **next Monday at 10h00**. However we may override this rule. In any cases, the dates indicated in Moodle are the one you should follow.

Exercise 2 Python installation on your computer

Skip this part if you are already set up with Python, Jupyter notebooks and your favorite IDE for Python. Otherwise, we recommend the following installations :

- Anaconda platform - a distribution of Python with popular data science packages that are pre-installed or easily installable. <https://www.continuum.io>. For the conservative party, select the 2.7 version of Python, for the daring party, select the 3.6 version of Python, as illustrated on Figure 1 below.

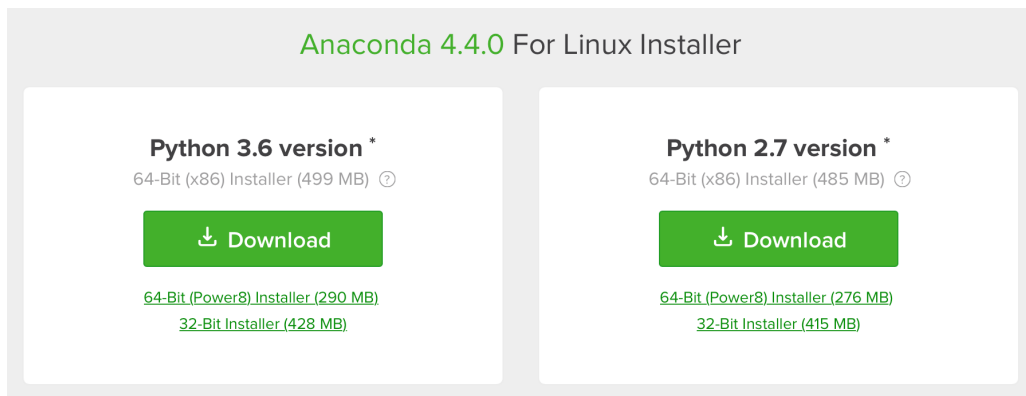


FIGURE 1 – Installing the Anaconda platform

- PyCharm IDE - an integrated development environment for Python. <https://www.jetbrains.com/pycharm>. You may select Professional or Community edition - the Community edition should be enough for this class. Free licenses are given to students, use your `hes-so.ch` address for the registration by JetBrains.

Exercise 3 Python language in a nutshell

We assume here that students are knowledgeable in other programming languages such as Java or C and that basic data structure concepts are known. If you know already Python and the concept of notebooks, then you can skip this exercise.

- Open Jupyter from the Anaconda Navigator. Jupyter Python notebooks run in your browser so, after launching Jupyter from Anaconda, a browser should show up. Open a new notebook from menu File and follow the User Interface Tour as illustrated in Figure 2.

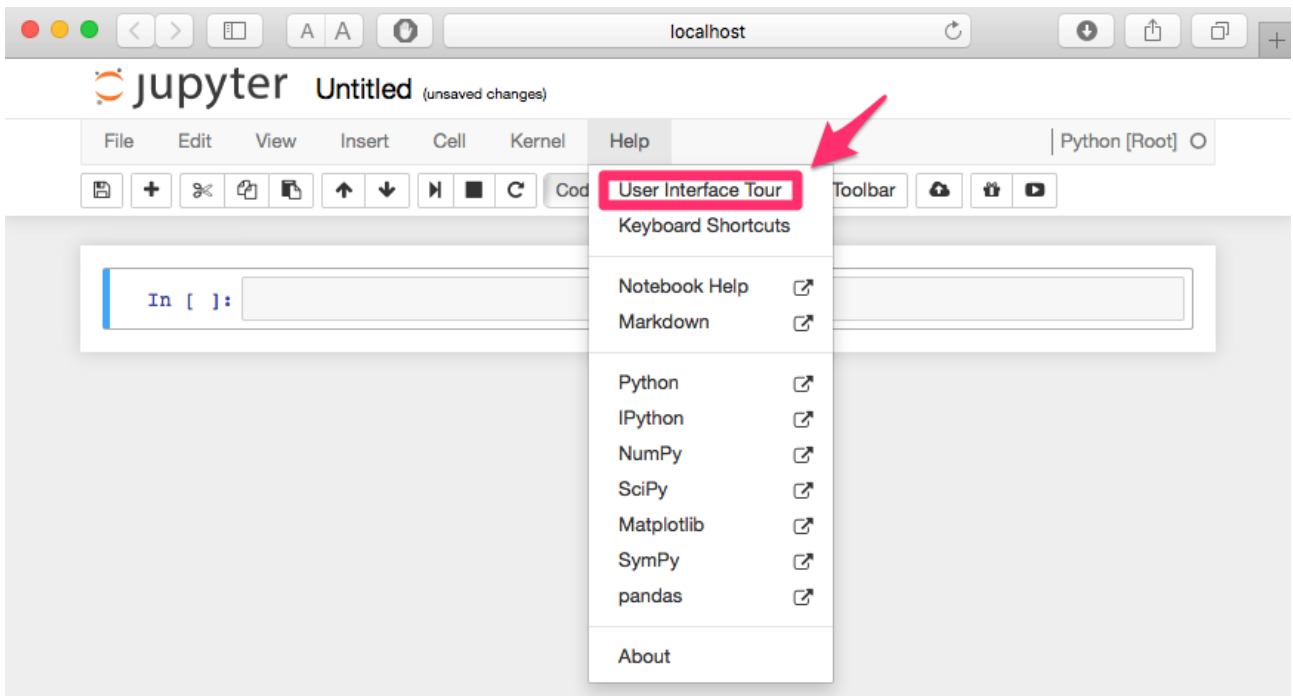


FIGURE 2 – First steps with Python notebooks.

- Download the file `intro-python.ipynb` from moodle and open it from Jupyter in Anaconda. You need to navigate where you stored the ipynb file. See Figure 3 below.

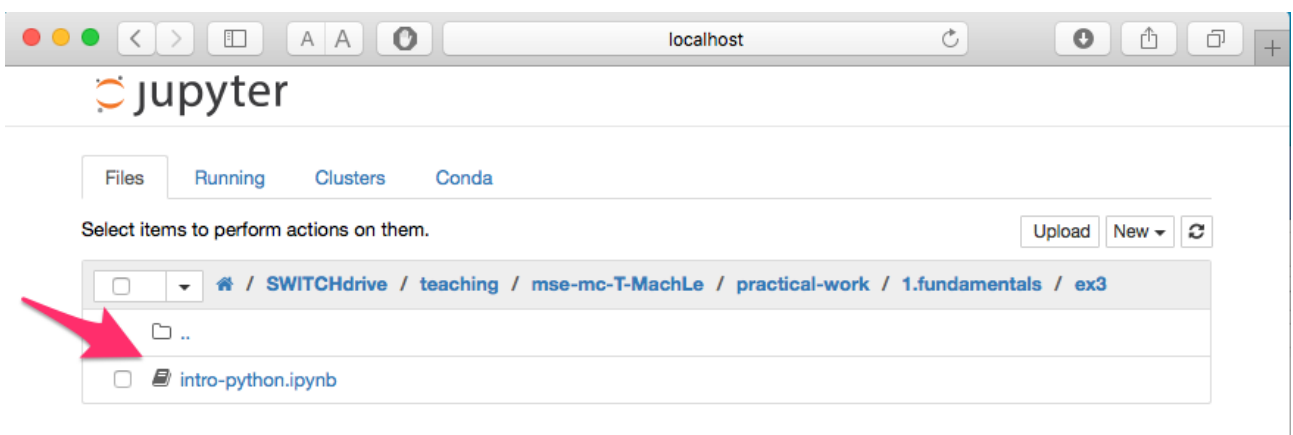


FIGURE 3 – Intro to Python language and Python notebooks.

- Go through the content of the `intro-python.ipynb` notebook and play with the cells. This document should give you a quick introduction to Python assuming that you are

fluent with other programming languages. You may also want to get familiar with [Mark-down syntax](#) if not known already.

- If you want a more fully fledged introduction to Python, read the official tutorial from <https://docs.python.org/2.7/tutorial/>.

Exercise 4 Data visualization

To train a bit yourself, we propose you to download the Iris dataset from <http://www.statlab.uni-heidelberg.de/data/iris/>, put it in a Python data structure and attempt to reproduce a plot close to the following one.

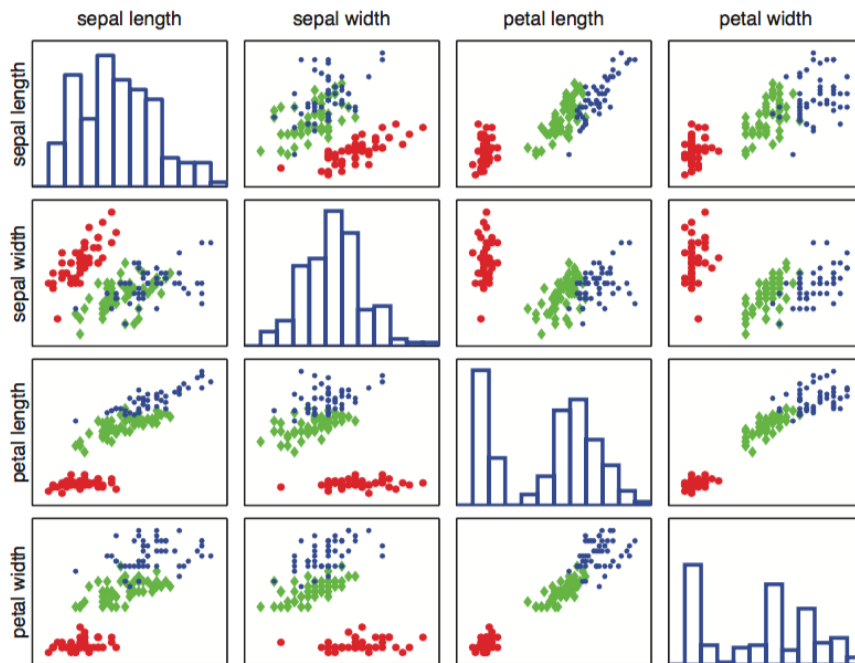


FIGURE 4 – Visualization of the Iris data as a pairwise scatter plot. The diagonal plots the marginal histograms of the 4 features. The off diagonals contain scatterplots of all possible pairs of features. Red circle = setosa, green diamond = versicolor, blue star = virginica.

The iris species classification task is a classical one. The goal is to distinguish three kinds of iris : setosa, versicolor and virginica. In the data set, a botanist has extracted 4 *features* : sepal length, sepal width, petal length and petal width.

What can you say regarding class separation? One class seems easier to distinguish from the others, which one? How would you separate the other two?

Exercice 5 Find your own examples of Machine Learning Tasks

Complete the Table below with 2 new examples of machine learning tasks, explaining what is the task T , performance measure P and training experience E .

Tableau 1-1. Quelques exemples de machine learning proposés par Mitchell

Cas d'application	Checkers learning	Handwriting recognition	Robot driving learning
Tasks T	Playing checkers	Recognizing and classifying handwritten words within images	Driving on public four-lane highways using vision sensors
Performance measure P	Percent of games won against opponents	Percent of words correctly classified	Average distance traveled before an error (as judged by human overseer)
Training experience E	Playing practice games againsts itself	A database of handwritten words with given classifications	A sequence of images and steering commands recorded while observing a human driver

FIGURE 5 – Source : Data Science - fondamentaux et études de cas, Michel Lutz et Eric Bernât, Eyrolles.

Exercice 6 Review questions

1) Supervised vs. unsupervised systems

Of the following examples, which one would you address using a supervised or an unsupervised learning algorithm ? Give some explanations for your answers.

- Given email labeled as spam/not spam, learn a **spam filter**.
- Given a set of news articles found on the web, group them into sets of **related articles**.
- Given a database of customer data, automatically discover **market segments** and group customers into different market segments.
- Given a dataset of patients diagnosed as either having **glaucoma** or not, learn to classify new patients as having glaucoma or not.

2) Classification vs. regression systems

Can we transform a regression problem into a classification problem ? What would be the benefits of doing so ?

Exercice 7 Optional - reading assignments

- Read the first chapter of Murphy’s book “Machine Learning”.
- Read the first chapter of Biernat and Lutz’s book “Data Science – fondamentaux et études de cas”.

Build your own summary of these chapters by doing a taxonomy of machine learning problems. You can find the pdfs on Moodle.