# Machine Learning with Python

## Numpy / Matplotlib / Scikit-learn

Giovanni Pellegrini

November 29, 2020

# Setup

## On lab machines



Download and extract the Scikit-learn lecture material from:

http://disi.unitn.it/~passerini/teaching/2020-2021/MachineLearning/

Open the terminal in the folder containing the extracted files and run:

```
> ./jupyter-scikit.sh
```

# Setup

Make sure you are using Python 3 for the following steps.

Install Numpy, Scipy, Matplotlib, Scikit-learn and Jupyter:

```
> pip install numpy scipy matplotlib sklearn
> pip install jupyter
```

Download and extract the material for the Scikit-learn lab:

http://disi.unitn.it/~passerini/teaching/2020-2021/MachineLearning/

Open the terminal in the folder containing the extracted files and run:

```
> jupyter notebook
```

# Setup
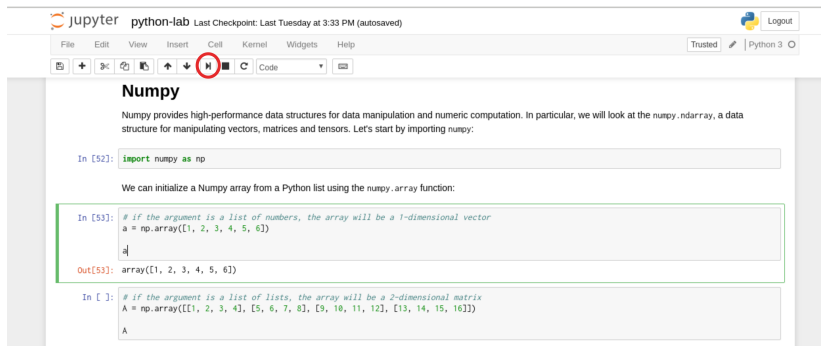
Open the browser at the given address and you'll see something like:



Open the `sklearn-lab.ipynb` file containing the lecture notebook.

# Jupyter notebook



Execute commands by selecting a cell and clicking the Run button on the header of the page or by **Shift+Enter**. You will see the output of the command just below the cell.

You can tweak and modify the code as you wish and execute it again.

# Exercise

For the exercise, you will solve a classification task using **Scikit-learn** over some given dataset. Each available dataset is already split into training and test sets. Choose a dataset, train a classifier on the training set and predict the labels on the test set. Hopefully, your classifier will classify the examples in the test set with higher accuracy than the reference baseline for the chosen dataset.

# Exercise

Datasets

**OCR**
Optical Character Recognition

**Spambase**
Spam email classification





**Presidential campaign tweets**
Classification of tweets from D. Trump and H. Clinton

# Exercise

Download the material:

`http://disi.unitn.it/~passerini/teaching/2020-2021/MachineLearning/`

The material contains the three datasets, each one containing:

▶ The training set examples;

▶ The training set labels;

▶ The test set examples;

▶ The test set labels;

▶ A README containing info about the dataset.
   this file also contains the reference baseline accuracy;

▶ Other info files.

# Exercise

1. Choose a dataset;
2. Experiment with a classification algorithm of your choosing;
3. Test your classifier using cross-validation over the training set
4. Train your classifier over the full training set;
5. Use the classifier to predict the examples in the test set;