# Introduction to Machine Learning for Chemists: Visualization, Data Processing, Analysis

**Romuald Poteau**
**Stella Christodoulou**
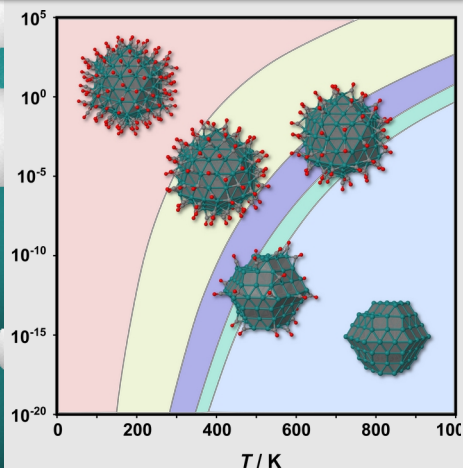
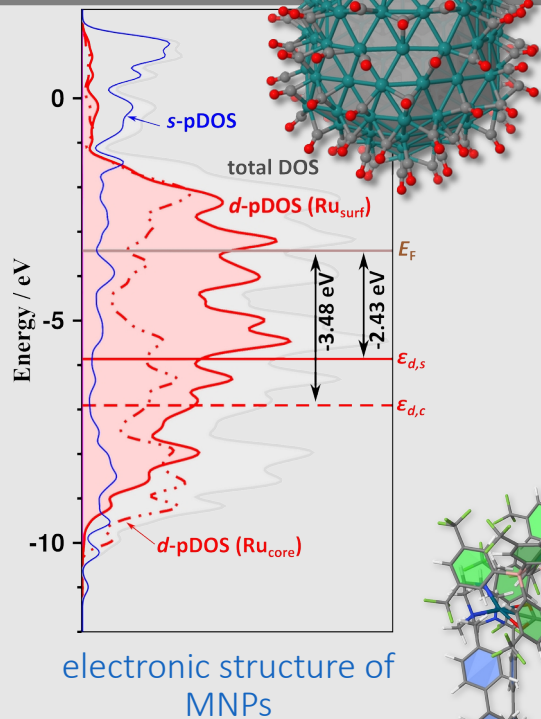## Applied Quantum Chemistry

DFT (Gaussian)
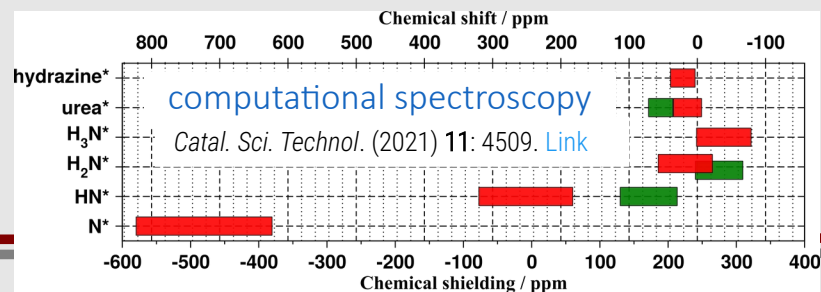DFT with PBC (VASP)
(basic tools freely available)
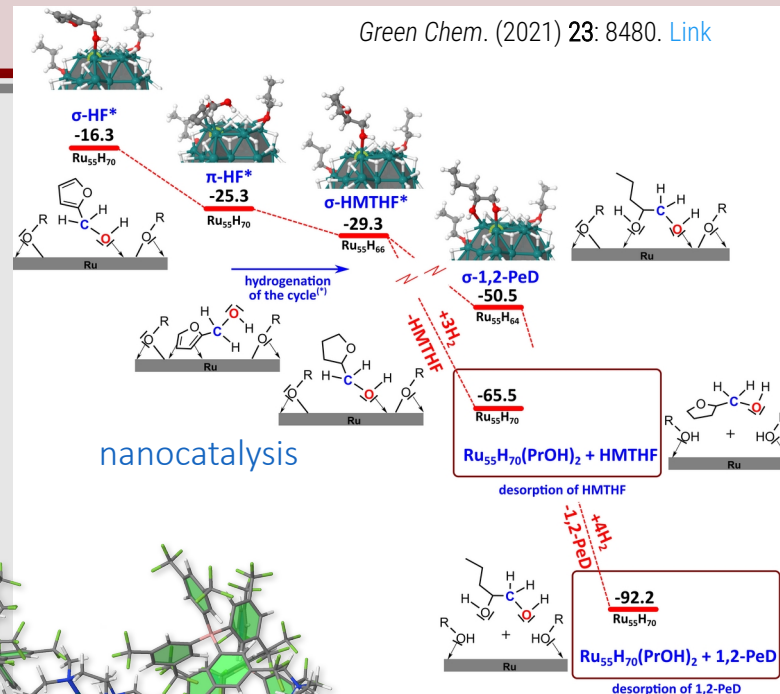


3D structure and stability of colloidal MNPs



electronic structure of MNPs

*Chem. Mater.* (2022) **34**: 2931-2944. Link

computational spectroscopy
*Catal. Sci. Technol.* (2021) **11**: 4509. Link



*Green Chem.* (2021) **23**: 8480. Link

nanocatalysis



*Nanoscale Horiz*. (2022) **7**: 607. Link

**"holistic" approach ← physical chemistry**

LPCNO

1. General introduction

2. Short selection of [simple] applications of supervised learning to chemistry

3. Tutorials / Live demonstrations ← Jupyter notebooks

github repository

Python in the
[Physical] Chemistry Lab

[ PytChem ]

https://github.com/rpoteau/PytChem

# General context

**Artificial Intelligence (AI)**
*intelligence demonstrated by machines, as opposed to the natural intelligence displayed by humans or animals*

**Goals**

reasoning & (basic) problem solving

knowledge representation

planning: making choices and hierarchy of events

learning (*i.e.* machine learning)

natural language processing

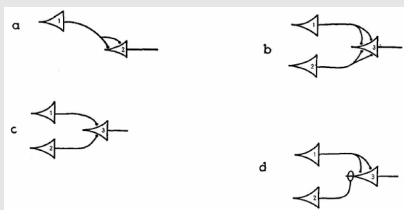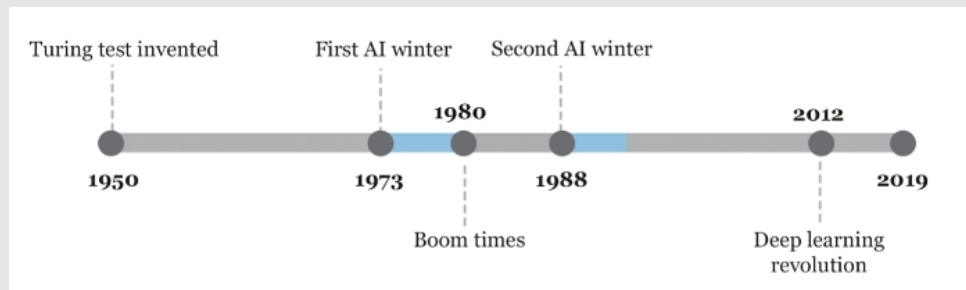perception of the world from sensors

ability to move and manipulate objects

simulating human affects


long-term goal: ability to solve an arbitrary problem


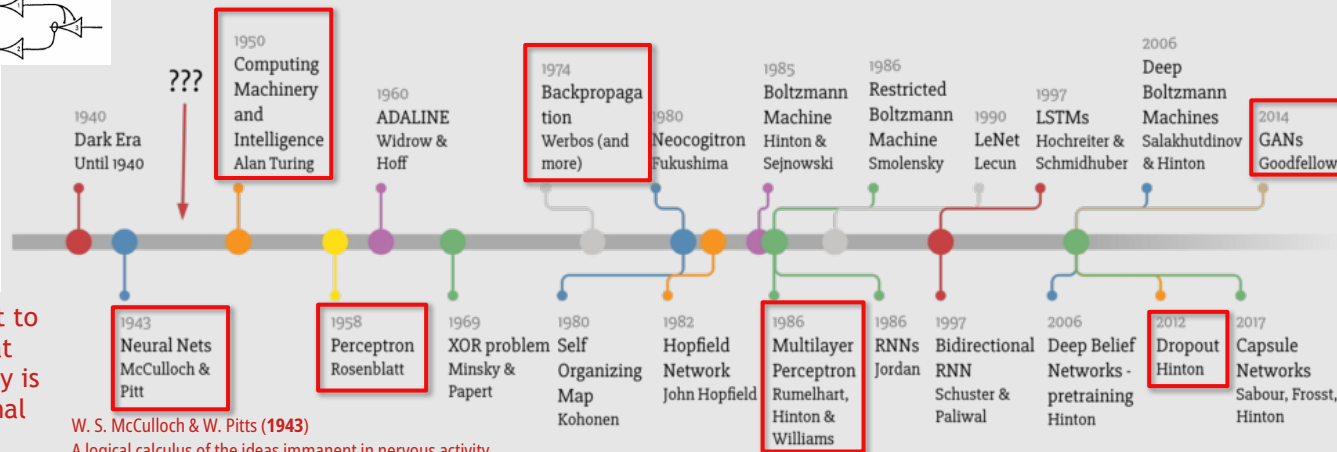*N.B. in some fields "artificial intelligence" means "machine learning with neural networks"*

Version: Sunday, October 1, 2023

LPCNO

## Artificial Intelligence & Deep learning timelines

**Since the first McC&P mathematical model for a neuron & the pioneering work of Turing, AI has a long history, with two "winters"**
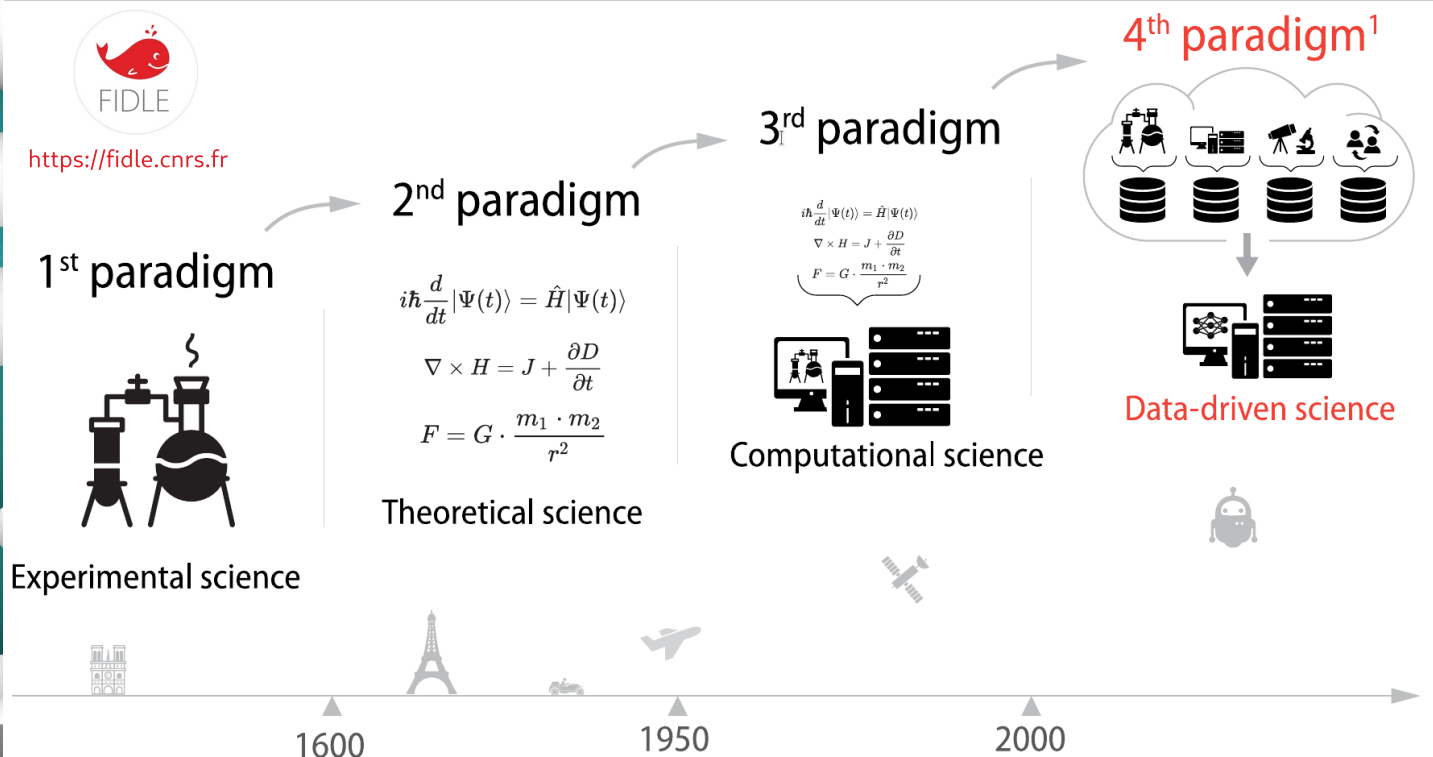


**neuron**

$x_1$
$x_2$
$x_3$
..
$x_n \in \{0,1\}$

$g$ $f$ $\longrightarrow y \in \{0,1\}$

positive or negative decision

were the first to suggest that neural activity is computational

W. S. McCulloch & W. Pitts (**1943**)
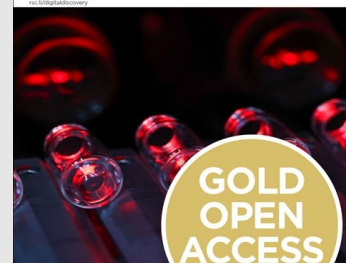A logical calculus of the ideas immanent in nervous activity
*Bull. Math. Biophys*. **5**, 115-133

Two neural networks contest with each other

**Version: Sunday, October 1, 2023**

## Artificial Intelligence and Scientific Research

NEW!



https://fidle.cnrs.fr

**1st paradigm**

Experimental science

**2nd paradigm**

$$i\hbar \frac{d}{dt}|\Psi(t)\rangle = \hat{H}|\Psi(t)\rangle$$

$$\nabla \times H = J + \frac{\partial D}{\partial t}$$

$$F = G \cdot \frac{m_1 \cdot m_2}{r^2}$$

Theoretical science

**3rd paradigm**

$$i\hbar \frac{d}{dt}|\Psi(t)\rangle = \hat{H}|\Psi(t)\rangle$$
$$\nabla \times H = J + \frac{\partial D}{\partial t}$$
$$F = G \cdot \frac{m_1 \cdot m_2}{r^2}$$

Computational science

**4th paradigm[1]**

Data-driven science

1600      1950      2000

**Digital Discovery**

GOLD OPEN ACCESS

"a new forum for data-driven approaches to scientific discoveries"

experimental and computational work

all topics related to the acceleration of discovery (screening, robotics, databases and advanced data analytics)

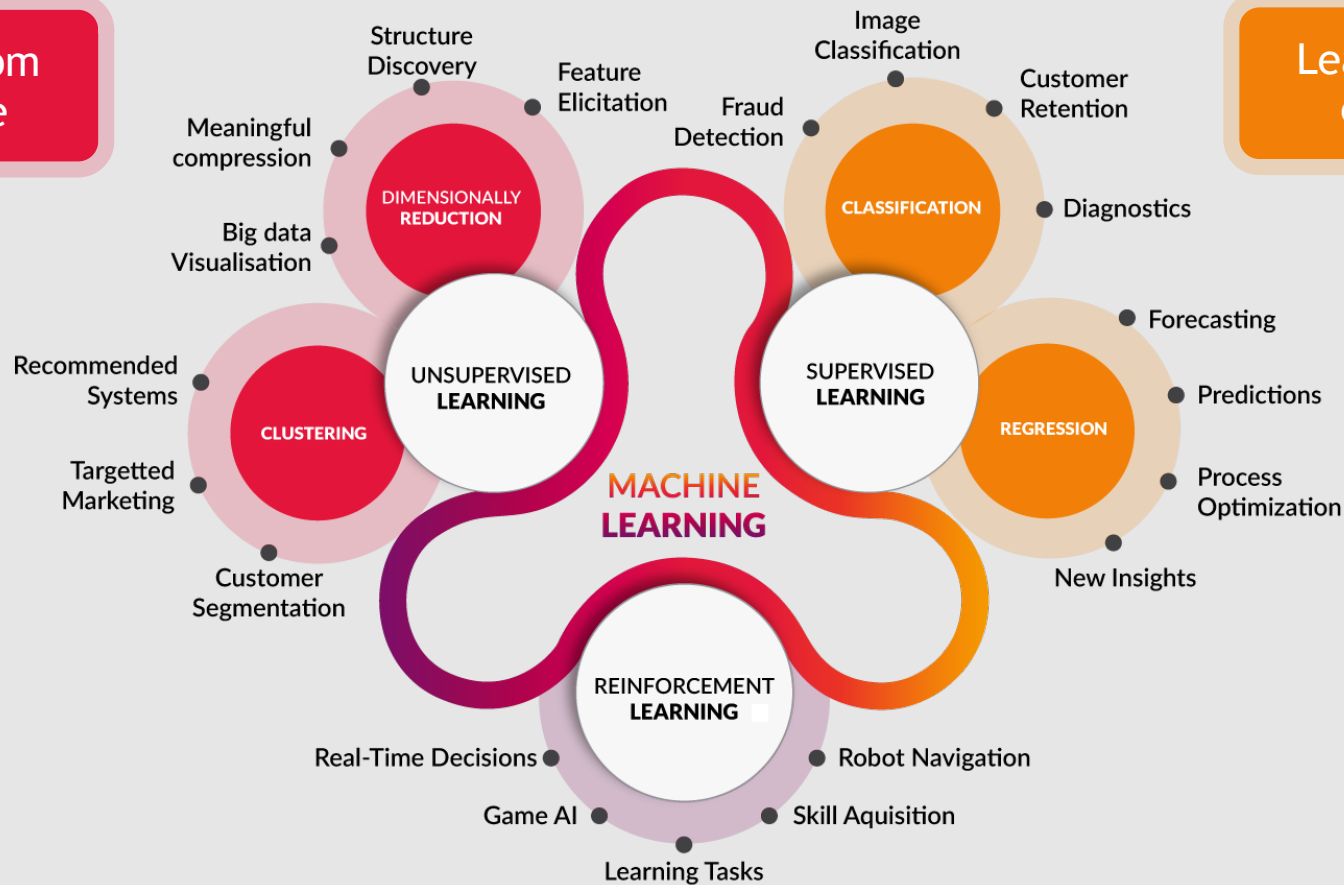broadly defined, but anchored in chemistry

(1) Hey, T.; Tansley, S.; Tolle, K. The Fourth Paradigm: Data-Intensive Scientific Discovery; The Fourth Paradigm: Data-Intensive Scientific Discovery; Microsoft Research, 2009
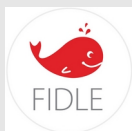
# Machine Learning

https://fidle.cnrs.fr

# Machine learning

**Data are provided along with the desired output (*i.e.* labelled data)**
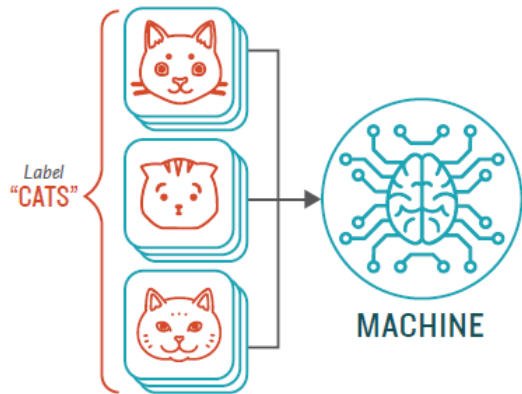
**Example of cats detection:**

    **- collect thousands of images of cats**

    **- draw a bounding box around each cat**

    **- feed the entire dataset to the machine so it can learn all by itself**
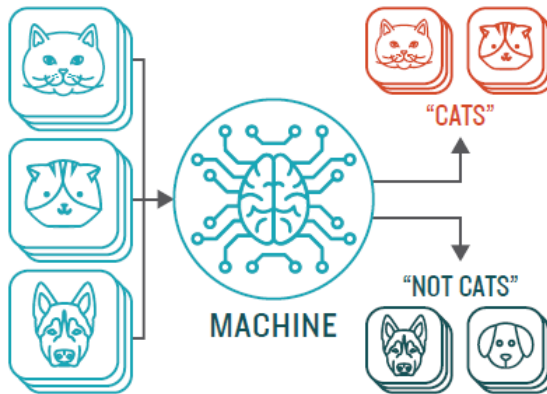
> **Learning from examples**



STEP 1 — Provide the machine learning algorithm categorized or "labeled" input and output data from to learn

STEP 2 — Feed the machine new, unlabeled information to see if it tags new data appropriately. If not, continue refining the algorithm
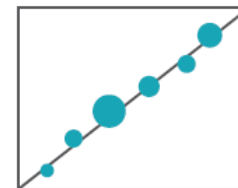
Label "CATS" → MACHINE

"CATS" / "NOT CATS" → MACHINE

**TYPES OF PROBLEMS TO WHICH IT'S SUITED**

**CLASSIFICATION** — Sorting items into categories

**REGRESSION** — Identifying real values (dollars, weight, etc.)

Version: Sunday, October 1, 2023

## Unsupervised learning

**Learning from data alone**

- **Just provide data**
- **Let the machine find out (or cluster) the patterns in the dataset**

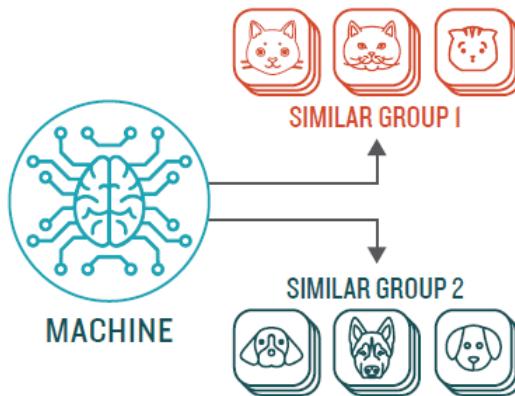# How to develop home-made ML tools?

mathematica

**Python is a high-level, interpreted, object-oriented, general-purpose programming language**

Core philosophy:

- Beautiful is better than ugly.
- Explicit is better than implicit.
- Simple is better than complex.
- Complex is better than complicated.
- Readability counts.

**~ 250 additional libraries are available**
- **data science**
- **machine learning**
- **modern scientific computation**
- **visualization**

**Using Python as an everyday tool for scientific calculation and computation**

**(it can even replace excel... by far)**

→ **basic knowledge of programming languages (variables, arrays, loops, conditional tests...)**

→ **enthusiastic and vast community**

→ **cheatsheets & cut/paste**

**Uneasy?**
- **yes and no**
- **important initial investment**
- **worth the effort**

python™



scikit **learn**

https://scikit-learn.org/
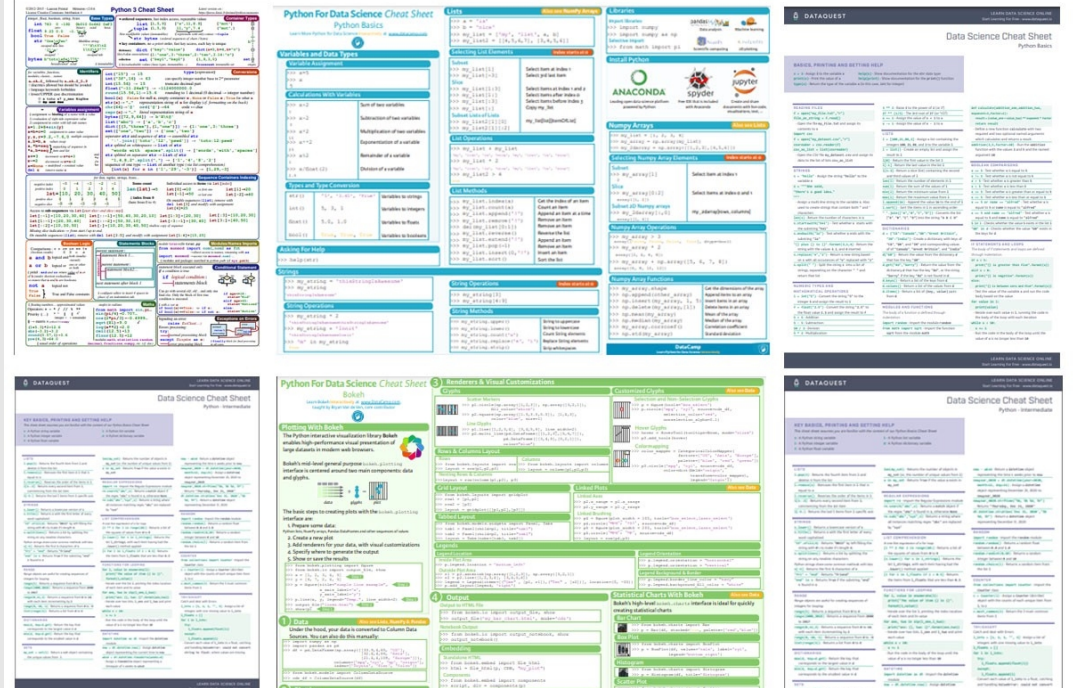
FUN.MOOC

Machine learning in Python with scikit-learn

## Simple and efficient tools for predictive data analysis
## Accessible to everybody, and reusable in various contexts
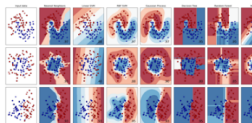## Built on NumPy, SciPy, and matplotlib

INRIA took leadership of the project and made the first public release on February 2010
3-clause BSD License (permissive free software license, compatible with the GNU GPL)

### Classification
Identifying which category an object belongs to.

**Applications:** Spam detection, image recognition.
**Algorithms:** SVM, nearest neighbors, random forest, and more...
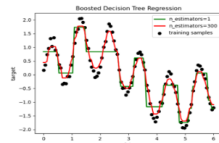
Examples

### Regression
Predicting a continuous-valued attribute associated with an object.

**Applications:** Drug response, Stock prices.
**Algorithms:** SVR, nearest neighbors, random forest, and more...

Examples

### Clustering
Automatic grouping of similar objects into sets.

**Applications:** Customer segmentation, Grouping experiment outcomes
**Algorithms:** k-Means, spectral clustering, mean-shift, and more...
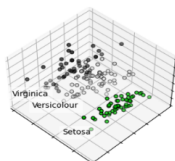
Examples

### Dimensionality reduction
Reducing the number of random variables to consider.

**Applications:** Visualization, Increased efficiency
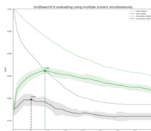**Algorithms:** k-Means, feature selection, non-negative matrix factorization, and more...

### Model selection
Comparing, validating and choosing parameters and models.

**Applications:** Improved accuracy via parameter tuning
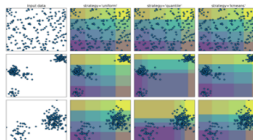**Algorithms:** grid search, cross validation, metrics, and more...

### Preprocessing
Feature extraction and normalization.

**Applications:** Transforming input data such as text for use with machine learning algorithms.
**Algorithms:** preprocessing, feature extraction, and more...

## K Keras
https://keras.io/

## High Level Deep Learning Application Programming Interface (API)
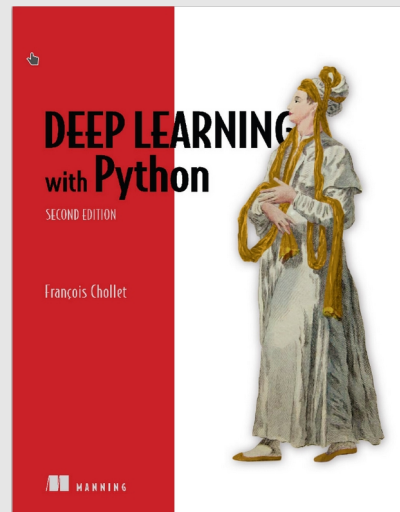
By François Cholet (Google)
Part on TensorFlow since 2017
MIT license (permissive free software license)

## how to start?

```
import numpy as np
import tensorflow as tf
from tensorflow import keras
```

### DEEP LEARNING with Python
SECOND EDITION
François Chollet
MANNING

TensorFlow
https://www.tensorflow.org/

## Google Brain's second-generation system

Supported by Google
Low level API
Apache license (yet another permissive free software license)

## PyTorch

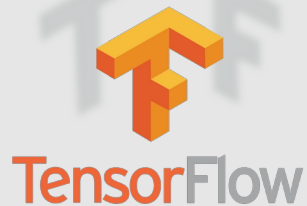https://pytorch.org/

From Torch library
Supported by Facebook
BSD licence
(permissive free software license)

**An open source machine learning framework that accelerates the path from research prototyping to production deployment**

**Widely used in the field of AI research**

ANACONDA

**ANACONDA DISTRIBUTION**

The world's most popular open-source Python distribution platform



Anaconda is a distribution of the Python and R programming languages for scientific computing (data science, machine learning applications, large-scale data processing, predictive analytics, etc.), that aims to simplify package management and deployment

JupyterLab is a web-based interactive development environment for notebooks, code, and data

Users can arrange workflows in **data science**, **scientific computing**, and **machine learning**

**you can save your everyday data manipulation / visualization as you do in your chemistry laboratory notebooks**

Version: Sunday, October 1, 2023