# Introduction to machine learning & Chemistry, or how I learned to be trendy

Clabaut Paul, Jiang Tao, Staub Ruben

ENS de Lyon

November 29,2019

# Theory & Generalities

LABORATOIRE
DE CHIMIE
ENS DE LYON

# What is machine learning ?

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T, as measured by P, improves with experience E."
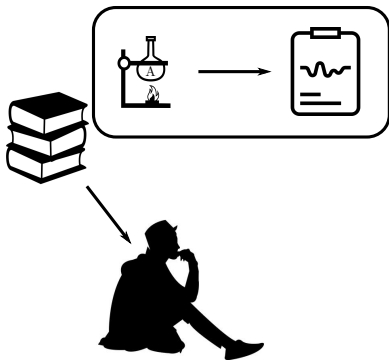
-Tom M. Mitchell

UNIVERSITÉ
DE LYON

Université Claude Bernard Lyon 1

cnrs

ENS DE LYON

# What is machine learning ?
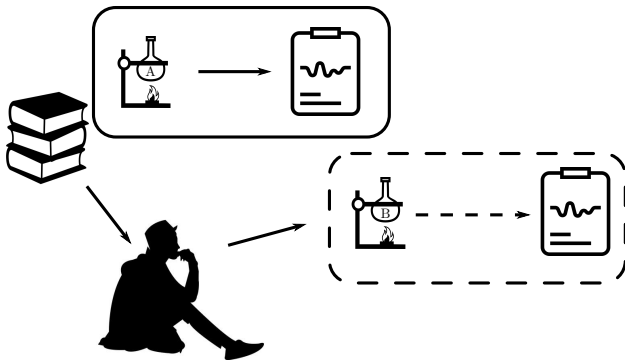
"Machine learning is glorified statistics"

"Machine learning is for Computer Science majors who couldn't pass a Statistics course."

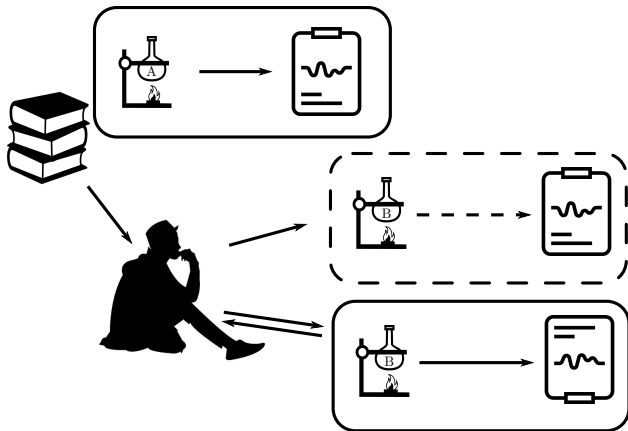"Machine learning is Statistics minus any checking of models and assumptions."
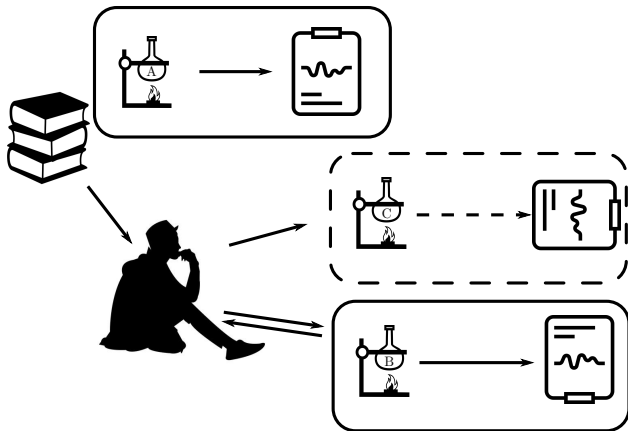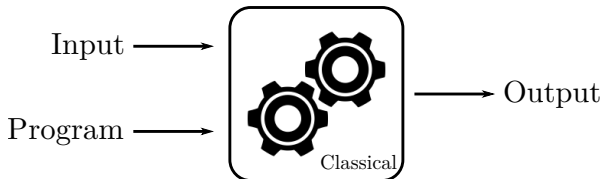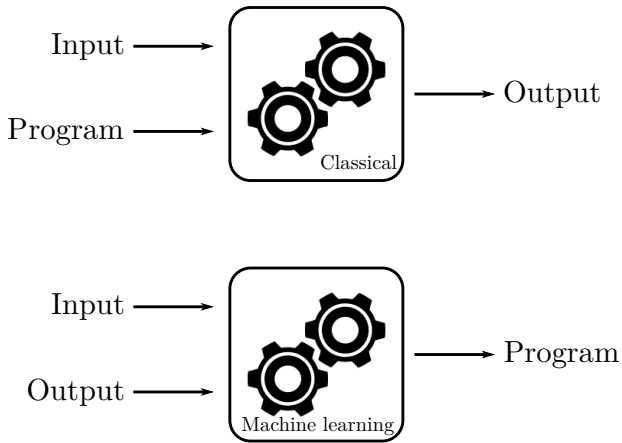
# What is machine learning ?

# What is machine learning ?

# What is machine learning ?

# What is machine learning ?

# What is machine learning ?
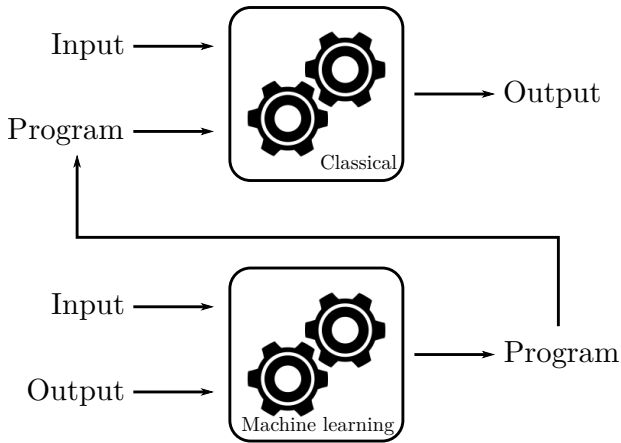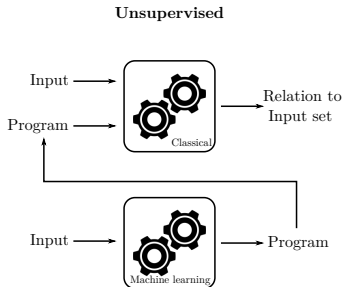
LABORATOIRE
DE CHIMIE
ENS DE LYON

# What is machine learning ?

# What is machine learning ?

# Supervised VS. Unsupervised

# Supervised VS. Unsupervised

# Families of Supervised learning

Aim: finding a function (program) mapping features (input) to a target (output)

Linguistic distinction based on target type:

Continuous target $\Rightarrow$ **Regression** (e.g. $E_{ads}$)

Categorical target $\Rightarrow$ **Classification** (e.g. active or not, cat or dog?)

Choose type of mapping function

# We all have done it a little before

# We all have done it a little before

# Linear model (regression)

Input: $m$ features, 1 target

|  | Features |  | Target |
|---|---|---|---|
|  | $(x_1 \ x_2 \ \cdots \ x_m)$ | $\rightarrow$ | $y$ |

Training set

| $(x_{11} \ x_{12} \ \cdots \ x_{1m})$ | $\rightarrow$ | $y_1$ |
|---|---|---|
| $(x_{21} \ x_{22} \ \cdots \ x_{2m})$ | $\rightarrow$ | $y_2$ |
| $\vdots$ |  |  |
| $(x_{n1} \ x_{n2} \ \cdots \ x_{nm})$ | $\rightarrow$ | $y_n$ |

Model: $m$ parameters $w_1, w_2, \ldots, w_m$

$$\sum_i^m w_i x_i \approx y$$

Pros: easily to understand

Cons: restricted to linear relations

LABORATOIRE
DE CHIMIE
ENS DE LYON

# Neural network (regression)

Simplest neural network

## Perceptron

Features

Output



$x_1$

$w_1$

$x_2$

$w_2$

$\vdots$

$x_m$

$w_m$

$b$

$f\left(\sum_i^m w_i x_i + b\right)$

Activation
function $f$

Parameters: $w_1, w_2, \ldots, w_m, b$

Equivalent to linear model if we use $f : x \mapsto x$

UNIVERSITÉ
DE LYON

Université Claude Bernard Lyon 1

CNRS

ENS DE LYON

LABORATOIRE
DE CHIMIE
ENS DE LYON

# Neural network (regression)

## Simplest neural network

### Perceptron

Features                                          Output



$f\left(\sum_i^m w_i x_i + b\right)$

$x_1$, $x_2$, ⋮, $x_m$ with weights $w_1$, $w_2$, $w_m$, bias $b$

Activation function $f$

## Common activation functions:

### Sigmoid

$\frac{1}{1+e^{-x}}$



### ReLU

$max(0, x)$

UNIVERSITÉ
DE LYON

Université Claude Bernard Lyon 1

CNRS

ENS DE LYON

LABORATOIRE
DE CHIMIE
ENS DE LYON

# Neural network (regression)

Input: $m$ features, $k$ targets



Parameters: weights $+$ bias of each neuron
Pros: fast to train, fast to predict, good fitting properties, trendy
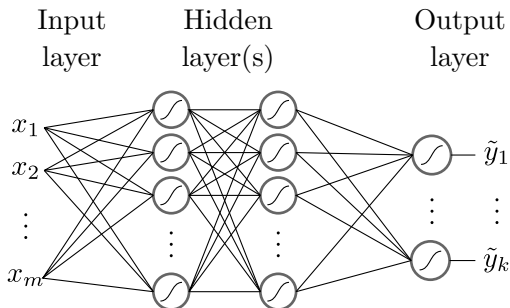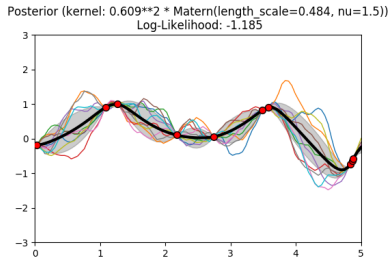Cons: black box, design tuning

UNIVERSITÉ
DE LYON

Université Claude Bernard Lyon 1    CNNS    ENS DE LYON

LABORATOIRE
DE CHIMIE
ENS DE LYON

# Gaussian process regression (regression)

Assume fitted target is a Gaussian process with given smoothness, use Bayesian inference to estimate probability distribution for prediction



Intuition: consider all possible fitted targets and extract probability distribution

Pros: native prediction confidence estimator

Cons: kernel-dependant, computationally costly training

UNIVERSITÉ
DE LYON

Université Claude Bernard Lyon 1

LABORATOIRE
DE CHIMIE
ENS DE LYON

# Support Vector Machine (classification)

Input: $m$ features, 1 categorical target



Intuition: Find best splitting hyperplane

UNIVERSITÉ
DE LYON

Université Claude Bernard    Lyon 1    cnrs    ENS DE LYON

LABORATOIRE
DE CHIMIE
ENS DE LYON

# Support Vector Machine (classification)

Change the metric through a non-linear kernel $\rightarrow$ embedding into higher-dimensional space where separation could be possible.



Pros: adapted for high-dimensionality
Cons: kernel-dependant, computationally costly training

UNIVERSITÉ
DE LYON

Université Claude Bernard Lyon 1

cnrs

ENS DE LYON

LABORATOIRE
DE CHIMIE
ENS DE LYON

## Decision tree (classification)

Input: *m* features, 1 categorical/continuous target



Intuition: Find best (feature+threshold) splitting on each node
Pros: easy interpretation (relevant features, ...), very fast prediction
Cons: requires balanced classes, optimal solution is NP-hard

UNIVERSITÉ
DE LYON

Université Claude Bernard Lyon 1    cnrs    ENS DE LYON

LABORATOIRE
DE CHIMIE
ENS DE LYON

# Decision tree (classification)

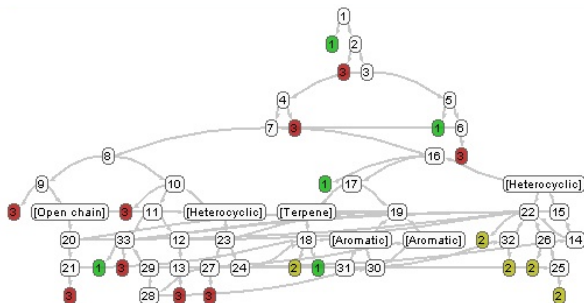Input: $m$ features, 1 categorical/continuous target



Intuition: Find best (feature+threshold) splitting on each node
Pros: easy interpretation (relevant features, ...), very fast prediction
Cons: requires balanced classes, optimal solution is NP-hard

UNIVERSITÉ
DE LYON

Université Claude Bernard Lyon 1      cnrs      ENS DE LYON

# Families of Unsupervised learning

Aim: Without consulting the target/output (sometimes absent), finding the structure in the input data itself.

Useful in certain cases:

* No prior knowledge of how many/what classes is the data divided into.

* Find out the most important features of the input data before feeding it to a machine.

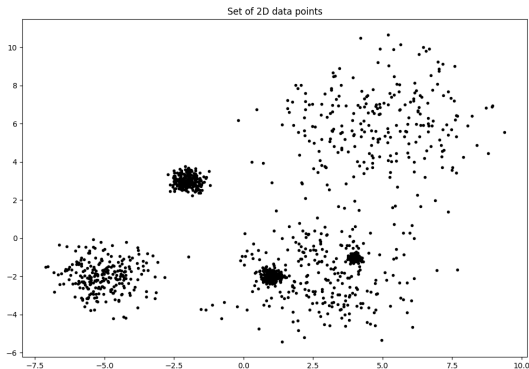Typical methods:

**Clustering**: k-means clustering, hierarchical clustering

**Dimensionality reduction**: PCA, Autoencoder

LABORATOIRE
DE CHIMIE
ENS DE LYON

# Clustering illustration

Clustering is partitioning into groups of close points
Non trivial task: How many clusters do you identify?



Set of 2D data points

UNIVERSITÉ
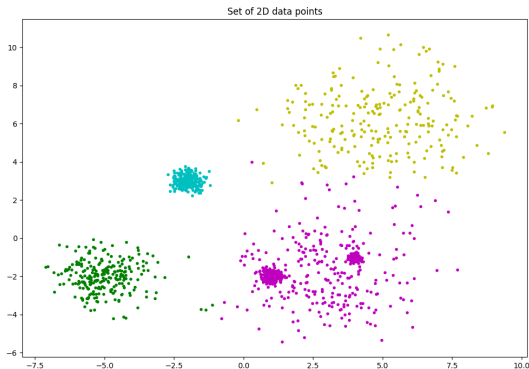DE LYON

Université Claude Bernard Lyon 1    cnrs    ENS DE LYON

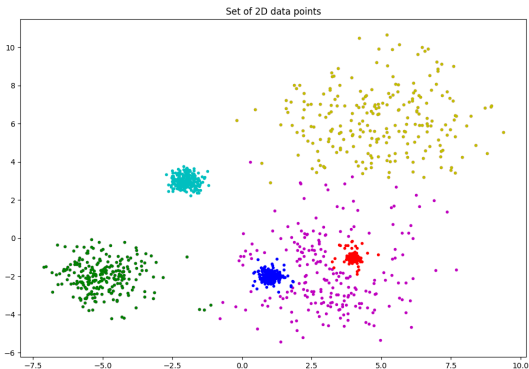# Clustering illustration

Clustering is partitioning into groups of close points
Non trivial task: How many clusters do you identify?



Set of 2D data points

# Clustering illustration

Clustering is partitioning into groups of close points
Non trivial task: How many clusters do you identify?



Set of 2D data points
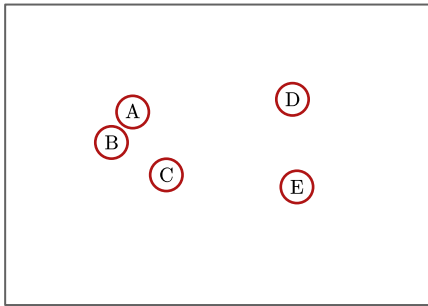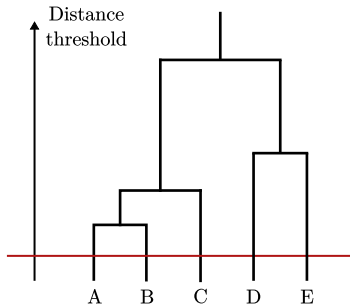
# Hierarchical clustering (clustering)

Input: metric between clusters, data points
Explore clusters generated for every threshold, produce dendrogram



Pros: easy to interpret, good overview
Cons: need user-defined metric between clusters, threshold selection

# Hierarchical clustering (clustering)

Input: metric between clusters, data points

Explore clusters generated for every threshold, produce dendrogram



Pros: easy to interpret, good overview

Cons: need user-defined metric between clusters, threshold selection

# Hierarchical clustering (clustering)

Input: metric between clusters, data points
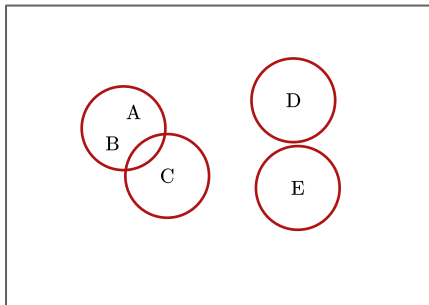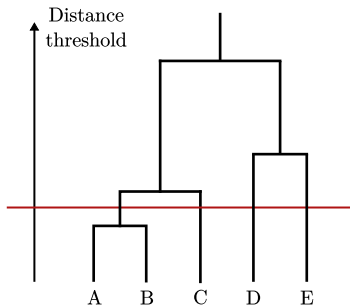Explore clusters generated for every threshold, produce dendrogram



Pros: easy to interpret, good overview
Cons: need user-defined metric between clusters, threshold selection

# Hierarchical clustering (clustering)

Input: metric between clusters, data points
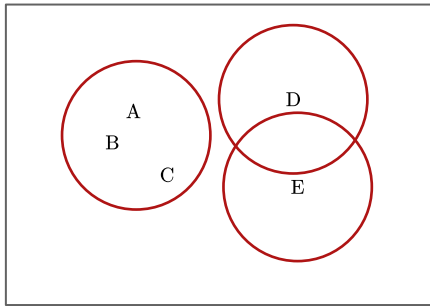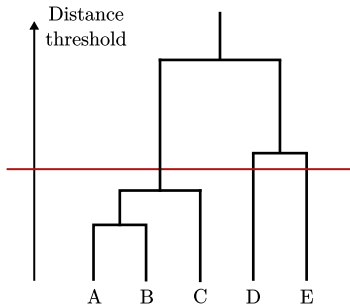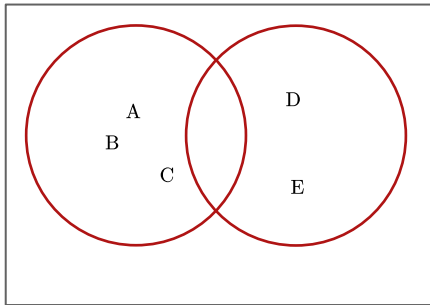Explore clusters generated for every threshold, produce dendrogram
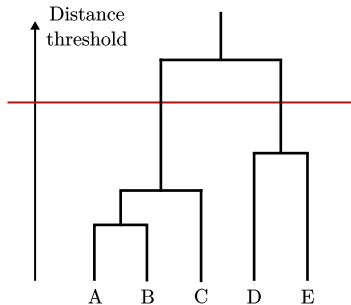


Pros: easy to interpret, good overview
Cons: need user-defined metric between clusters, threshold selection

# K-means clustering (clustering)

Input: number of clusters, distance between data points
Find centroids that minimize the within-cluster sum-of-squares



Pros: easy to understand
Cons: stochastic, assume convex and isotropic, poor
high-dimensionality support

# PCA (dimensionality reduction)



Emphasize variation and bring out strong patterns in a dataset.

Useful for finding important features in high dimensional dataset.

# PCA (Wine chemistry)

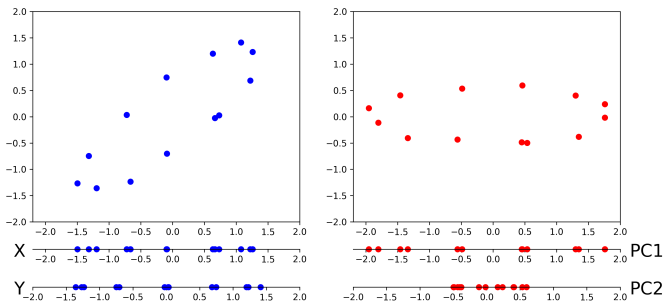| Feature No. | Composition |
|---|---|
| 1 | Alcohol |
| 2 | Malic acid |
| 3 | Ash |
| 4 | Alcalinity of ash |
| 5 | Magnesium |
| 6 | Total phenols |
| 7 | Flavanoids |
| 8 | Nonavanoid phenols |
| 9 | Proanthocyanins |
| 10 | Color intensity |
| 11 | Hue |
| 12 | OD280/OD315 |
| 13 | Proline |

| Class | Number of wines |
|---|---|
| 1 | 59 |
| 2 | 71 |
| 3 | 48 |

https://archive.ics.uci.edu/ml/datasets/Wine, Accessed: 2019-11-29

LABORATOIRE
DE CHIMIE
ENS DE LYON

# PCA (Wine chemistry)

| Feature No. | Composition |
|---|---|
| 1 | Alcohol |
| 2 | Malic acid |
| 3 | Ash |
| 4 | Alcalinity of ash |
| 5 | Magnesium |
| 6 | Total phenols |
| 7 | Flavanoids |
| 8 | Nonavanoid phenols |
| 9 | Proanthocyanins |
| 10 | Color intensity |
| 11 | Hue |
| 12 | OD280/OD315 |
| 13 | Proline |

UNIVERSITÉ
DE LYON

Université Claude Bernard Lyon 1    cnrs    ENS DE LYON

# Manifold (dimensionality reduction)

Non-linear dimensionality reduction:

  Kernel PCA (use kernel instead of covariance)

  MDS (preserve distances)

  Isomap (preserve geodesic graph-based distances)
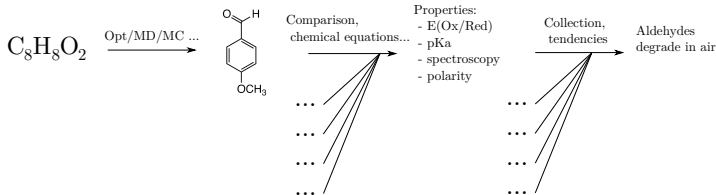
  Self-organizing maps (preserve topology)

  . . .

# Application

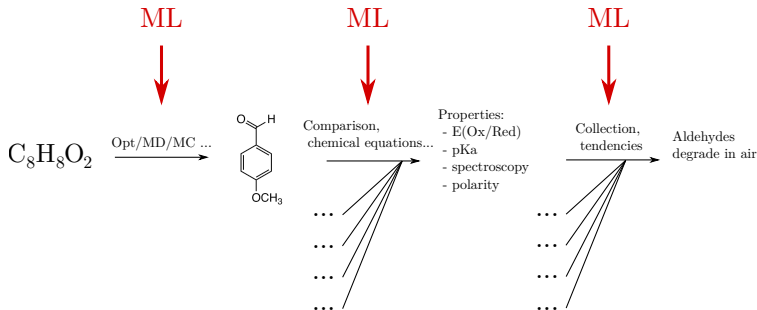LABORATOIRE
DE CHIMIE
ENS DE LYON

# Categories of applications in chemistry

Substitute to computation algorithm, Direct property prediction, Data analysis

# Categories of applications in chemistry

Substitute to computation algorithm, Direct property prediction, Data analysis

# Application: Water/Platinum potential: Goal

## DFT
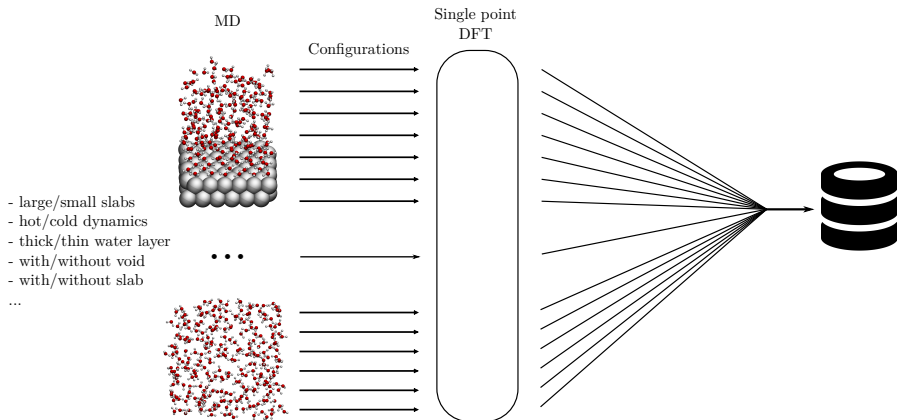
500 ps of 20Å thick water layer on 3*3*4 Pt (111) slab
$\rightarrow$ 8 yrs on 100 processors

## MM

$E = E_{slab} + \sum\limits_{wat}(E_{wat} + E_{slab/wat}) + \sum\limits_{wat}\sum\limits_{wat}(E_{wat/wat} + E_{slab/wat/wat}) + \sum\limits_{wat}\sum\limits_{wat}\sum\limits_{wat}(...) + ...$
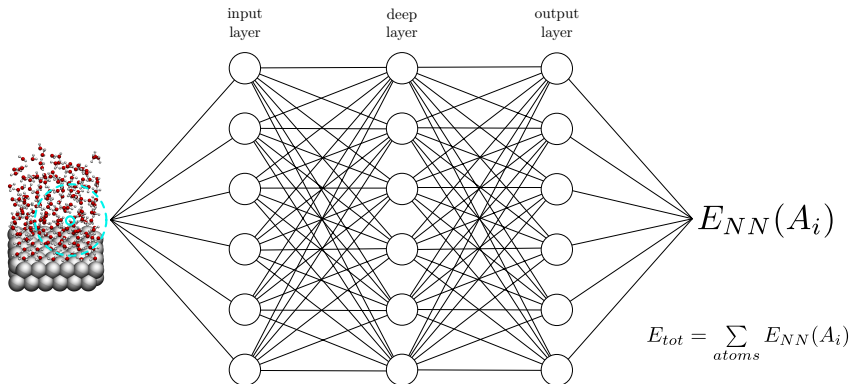
## Neural Network

$E = \sum\limits_{atoms} E_{NN}(environement) \rightarrow$ Might just work !

Université Claude Bernard Lyon 1

# Step 1 Gathering data



MD

Configurations

Single point
DFT

- large/small slabs
- hot/cold dynamics
- thick/thin water layer
- with/without void
- with/without slab
...

# Step 2 Input and Training the Neural Network



input
layer

deep
layer

output
layer

$$E_{NN}(A_i)$$

$$E_{tot} = \sum_{atoms} E_{NN}(A_i)$$

# Step 2 Input and Training the Neural Network



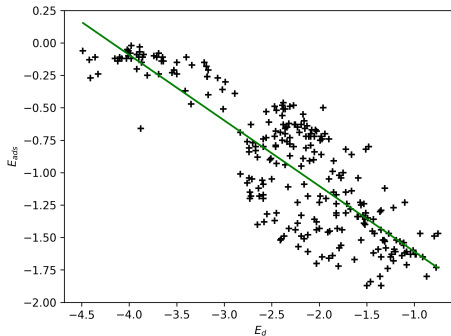input layer <=>
ensemble of symetry function

$F_{radial}(H_i, r, O) = 2$

$F_{radial}(H_i, r, H) = 4$
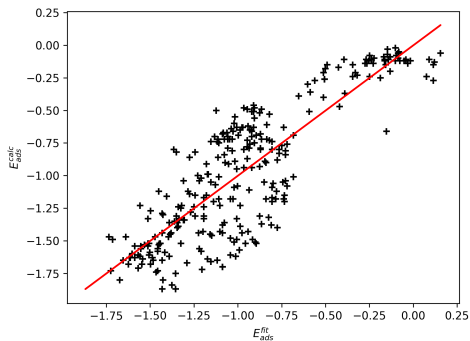
$N_{coefficient}$

$= (N_{input\_function} * N_{neuron,layer1}$

$+ N_{neuron,layer1} * N_{neuron,layer2}$

$+ N_{neuron,layer2}) * 3$ (O,H,Pt)

Artrith, N.; Behler, J. *Physical Review B* **2012**, *85*, DOI:
10.1103/PhysRevB.85.045439

LABORATOIRE
DE CHIMIE
ENS DE LYON

# The adsorption energy example

UNIVERSITÉ
DE LYON

Université Claude Bernard Lyon 1

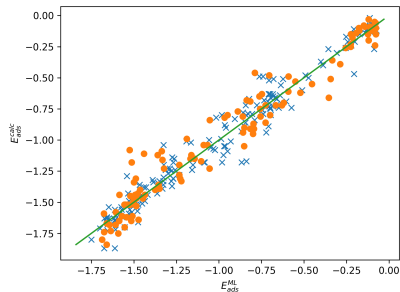CNRS

ENS DE LYON

# Linear fit vs calculation



Only one descriptor ($E_d$) is not good enough!
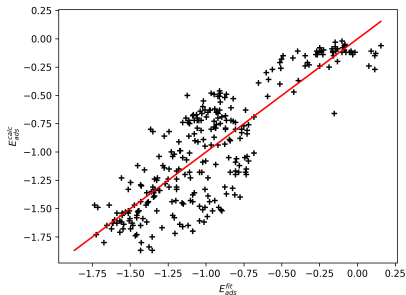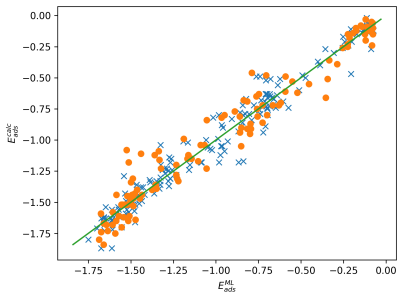Even if it is based on a physical model (tight binding model)

# What if we consider more descriptors?

| | | |
|---:|:---:|:---|
| $f$ | : | Filling of a d-band |
| $E_d$ | : | Center of a d-band |
| $W_d$ | : | Width of a d-band |
| $\gamma_1$ | : | Skewness of a d-band |
| $\gamma_2$ | : | Kurtosis of a d-band |
| $W$ | : | Work function |
| $r_0$ | : | Atomic radius |
| $r_d$ | : | Spatial extent of d-orbitals |
| IE | : | Ionization potential |
| EA | : | Electron affinity |
| $\chi_0$ | : | Pauling electronegativity |
| $\chi$ | : | Local Pauling electronegativity |
| $V_{ad}^2$ | : | Interatomic d coupling matrix element |

## ML vs calculation

Ma, X. et al. *The Journal of Physical Chemistry Letters* **2015**, *6*, 3528–3533

ML significantly improves the fit by utilizing many descriptors.

ML needs input with chemistry insight in it.

**ML is a tool, not magic.**

Ma, X. et al. *The Journal of Physical Chemistry Letters* **2015**, *6*, 3528–3533

# Domain of applicability

Models must be treated with care

Beware of **overfitting**:
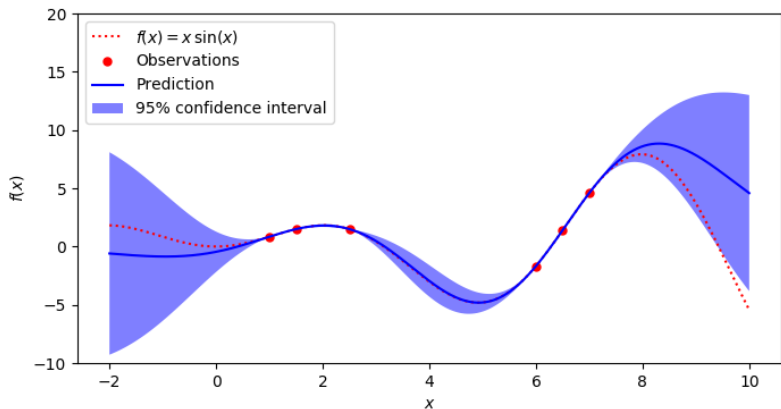
Especially with high-dimensionality

Quality estimator: **cross-validation** is a good starting point

Learned models are meant for **interpolation**, not extrapolation:

Would require physico-chemical justification

Training set should cover your subsequent usage

LABORATOIRE
DE CHIMIE
ENS DE LYON

# Extrapolation illustration

LABORATOIRE
DE CHIMIE
ENS DE LYON

# Checking assumptions
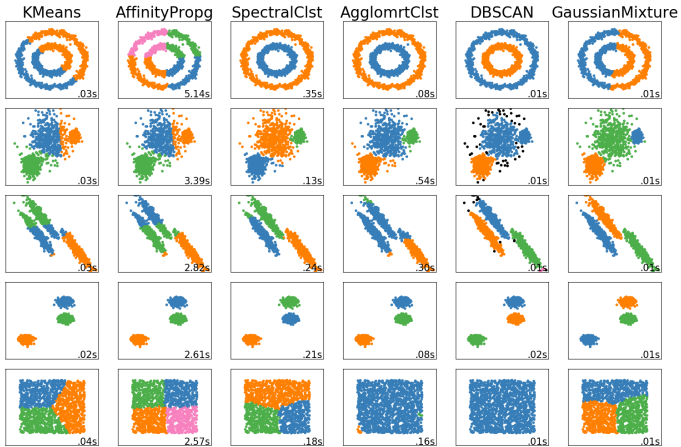
### Theorem (No free lunch)

*Any two optimization algorithms are equivalent when their performance is averaged across all possible problems*

$\Rightarrow$ There cannot exist a machine learning algorithm that outperforms all other algorithms on every problem
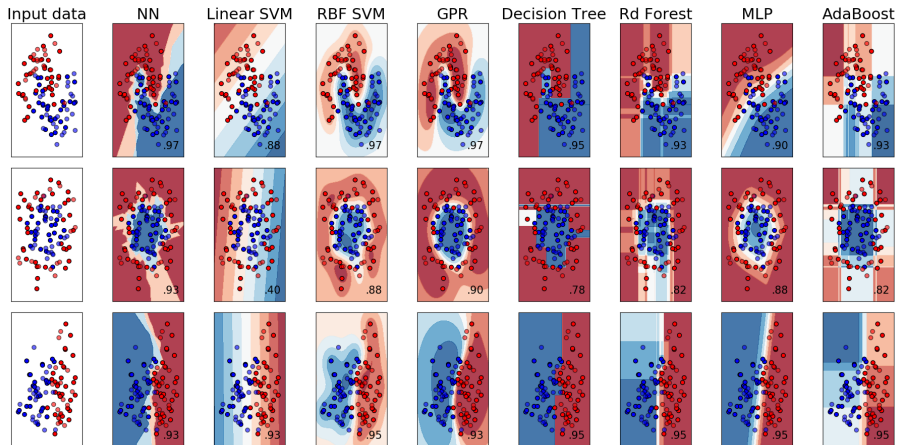
$\Rightarrow$ A machine learning algorithm can be better than an other **only** under specific assumptions

The moral being: compare and select the algorithm that fits the best your problem

# Clustering comparison

Classifiers comparison

# Another paradigm?

Third basic paradigm of machine learning: **reinforcement learning**

Aim: apply best policy to minimize regrets, without initial expertise (learn policies on the fly)

Trade-off between exploration and exploitation

Applications: decision making, global minimization, …

Algorithms: MCTS, genetic algorithms, …

# Grandmaster level in StarCraft II using multi-agent reinforcement learning.png

MENU ⌄

## nature

Article | Published: 30 October 2019

# Grandmaster level in StarCraft II using multi-agent reinforcement learning

Oriol Vinyals ✉, Igor Babuschkin, […] David Silver ✉

**36k** Accesses | **1** Citations | **916** Altmetric | Metrics

# The Singularity is Near