# *Initiation à l'apprentissage automatique en science des matériaux*
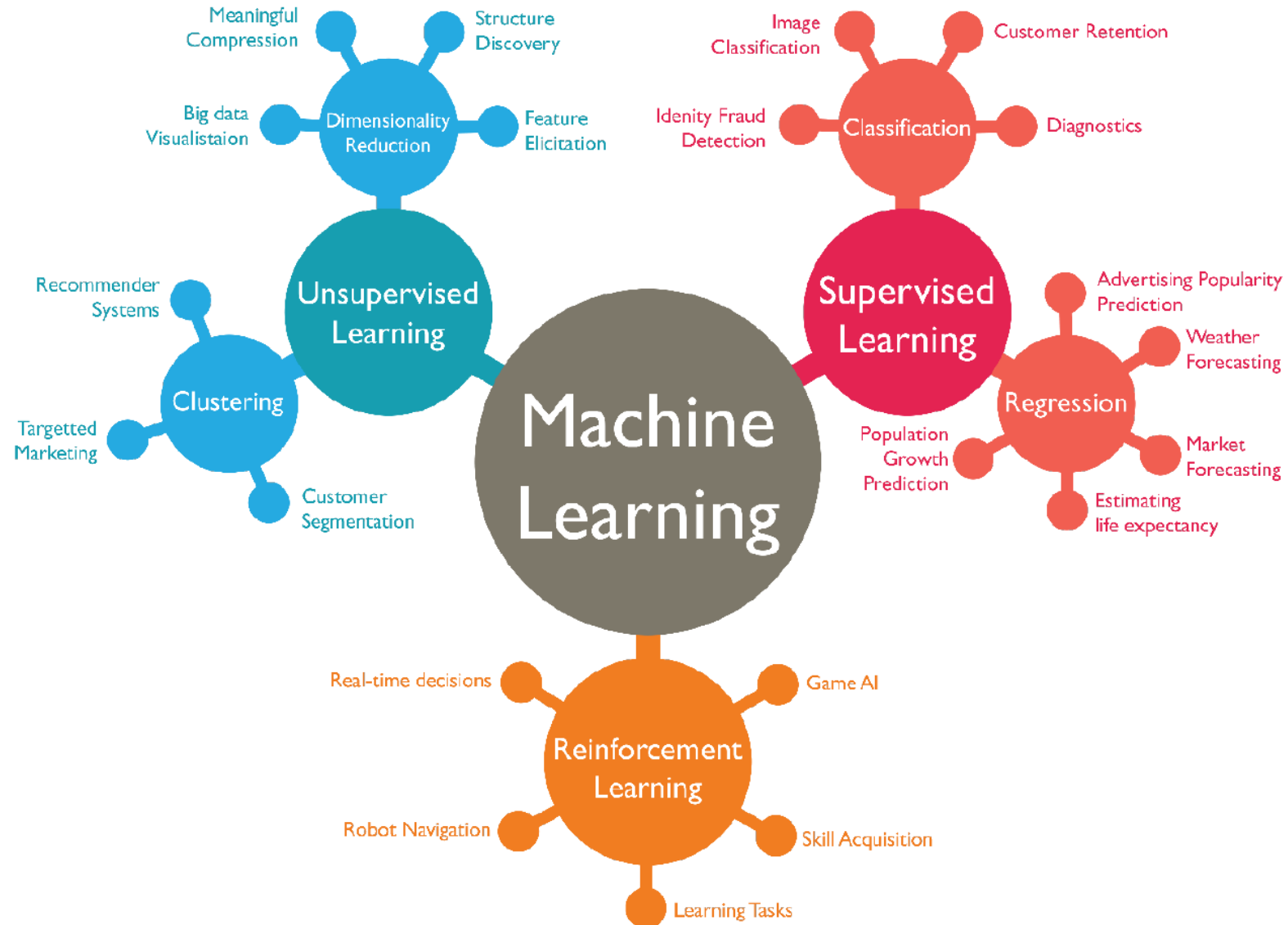
## 2. Fundamentals in Machine Learning

J.-C. Crivello, LINK : jean-claude.crivello@cnrs.fr
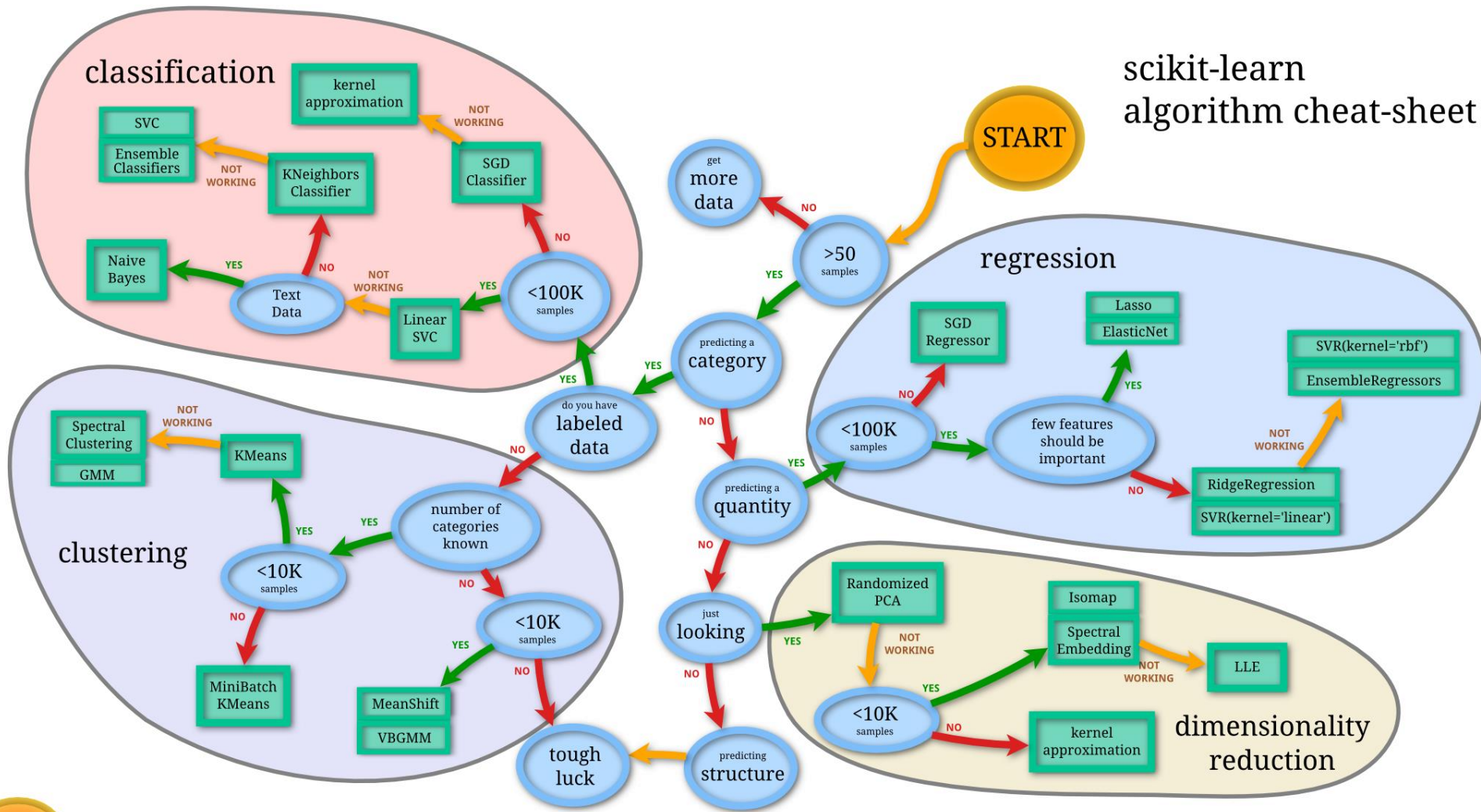
C. Barreteau, ICMPE : celine.barreteau@cnrs.fr

S. Junier, ICMPE : sebastien.junier@cnrs.fr

https://link.cnrs.fr/ML/
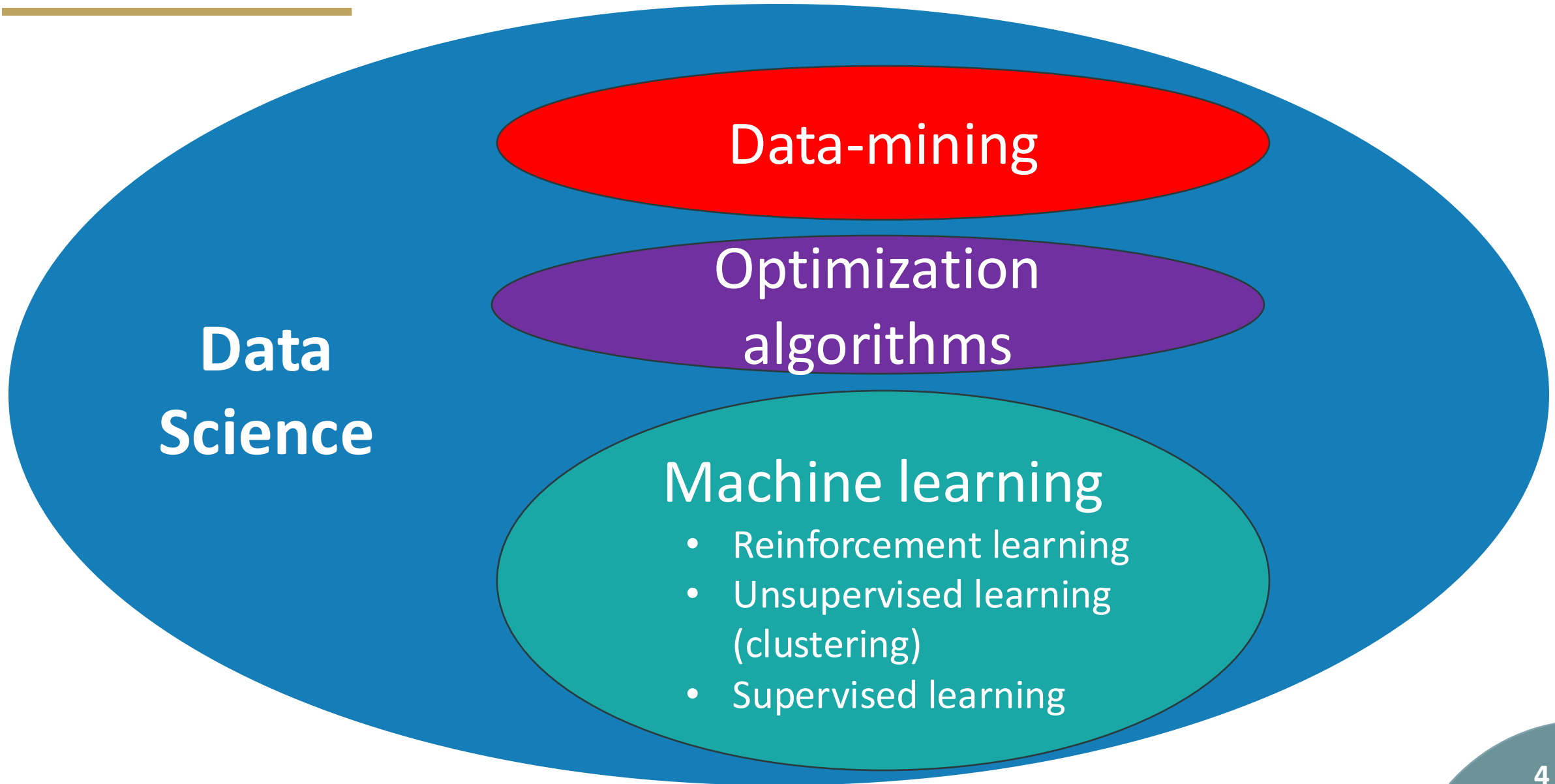
# Machine learning

# Machine learning



scikit-learn algorithm cheat-sheet

**Choosing the right estimator?**
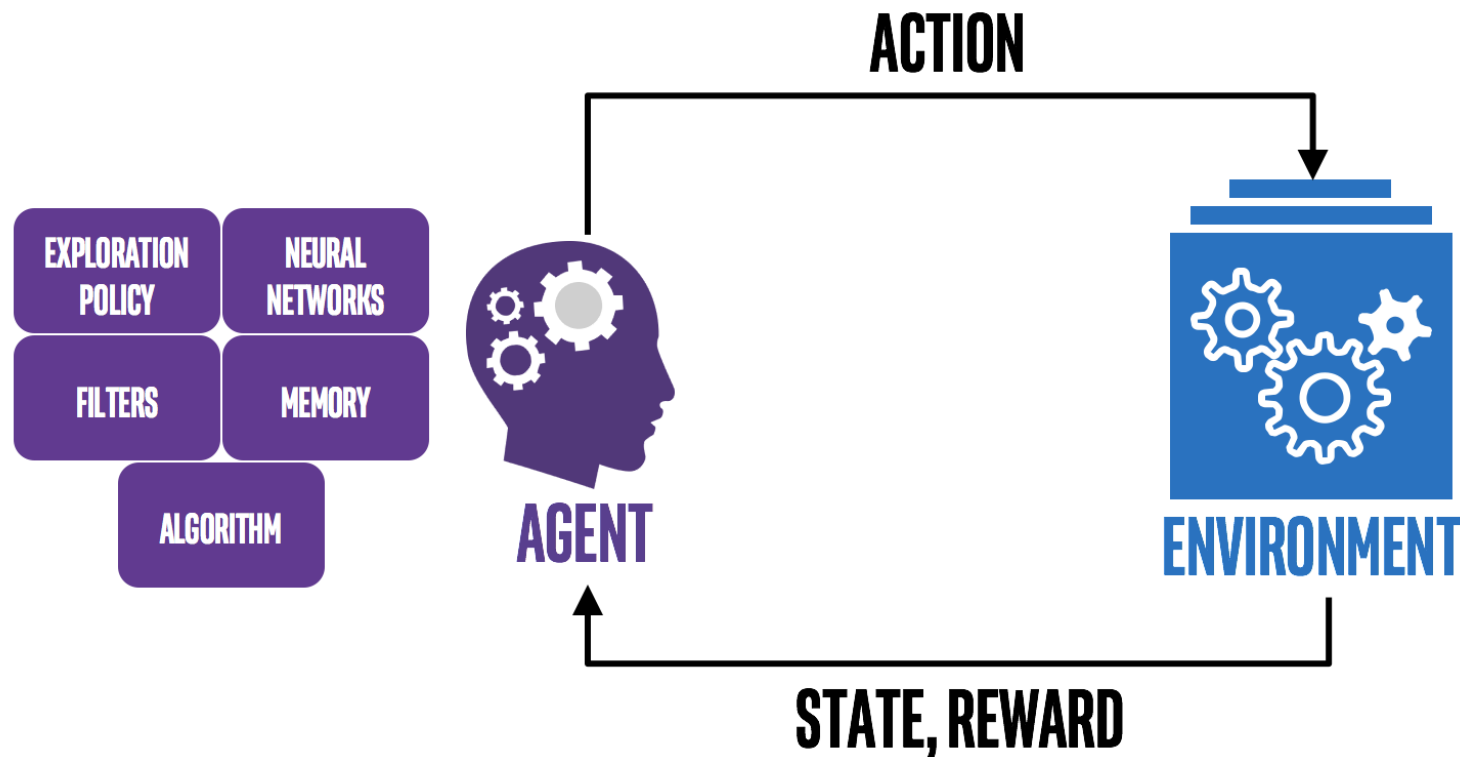https://scikit-learn.org/stable/tutorial/machine_learning_map/index.html

3

# Several approaches of the data science



Data Science

Data-mining

Optimization algorithms

Machine learning
- Reinforcement learning
- Unsupervised learning (clustering)
- Supervised learning

# 1. Reinforcement learning

RL is learning from experiences.

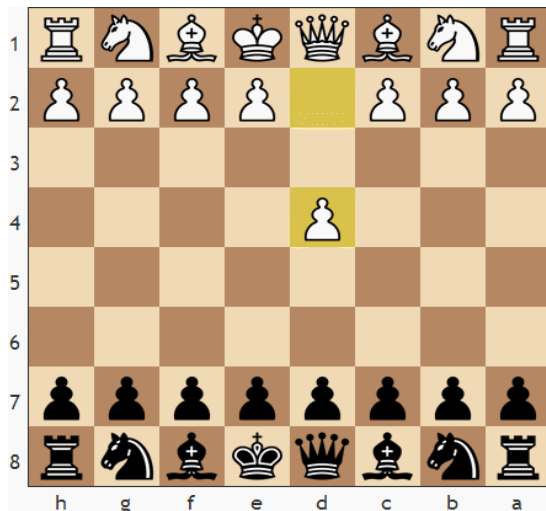RL teaches an agent how to choose an action from its action space, within a particular environment, in order to maximize rewards over time
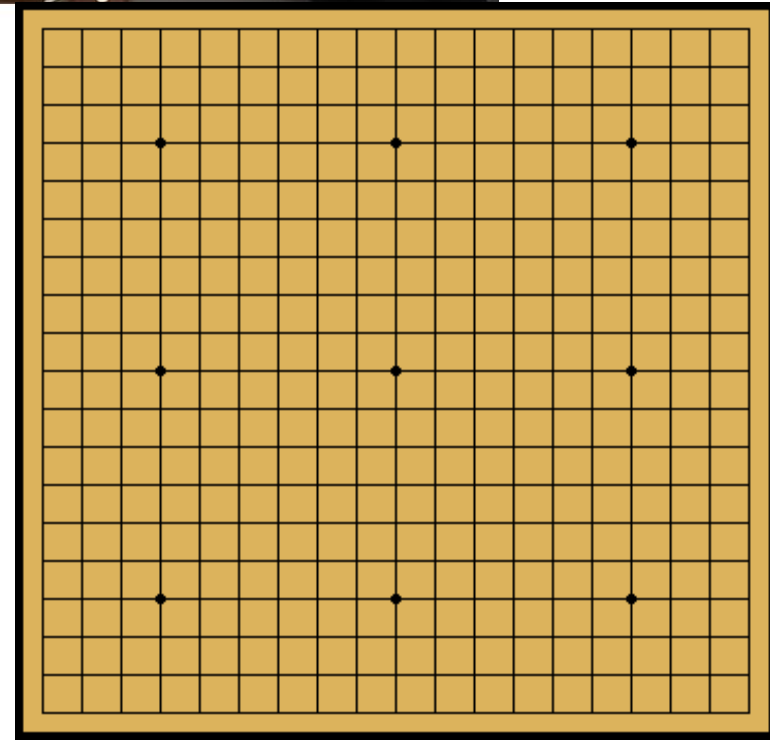
ACTION

EXPLORATION POLICY

NEURAL NETWORKS

FILTERS

MEMORY

ALGORITHM

AGENT

ENVIRONMENT

STATE, REWARD

# Game theory: Brute force VS RL


1997 – Deep Blue - IBM


Google DeepMind
Challenge Match
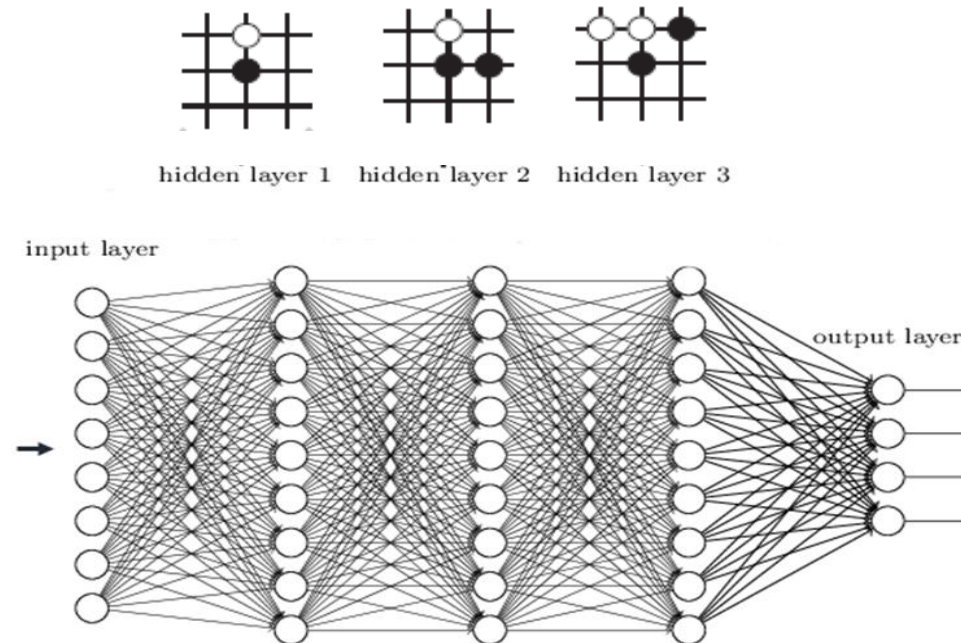8 - 15 March
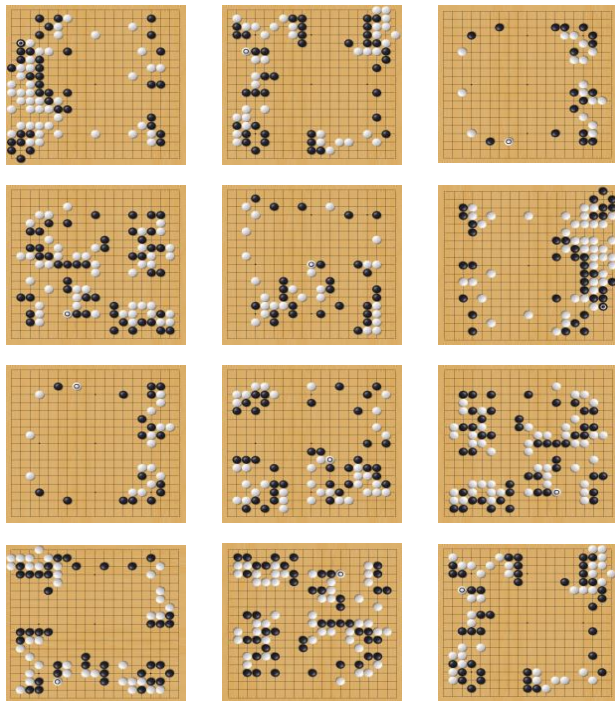2016
Alpha Go

**Number of possible sequences:**
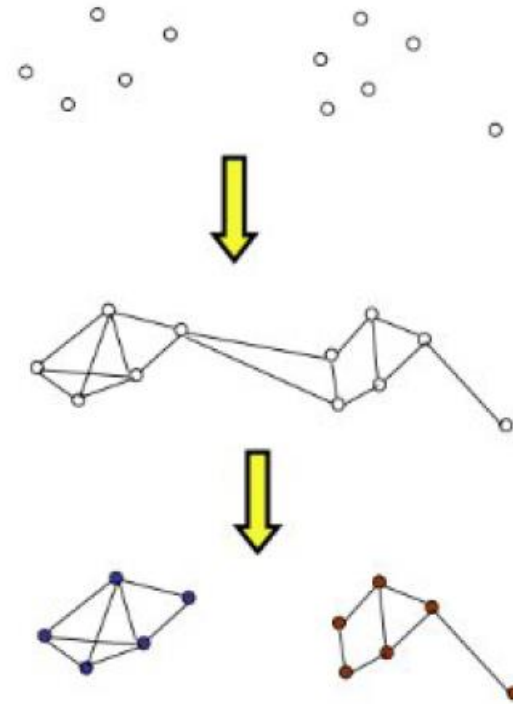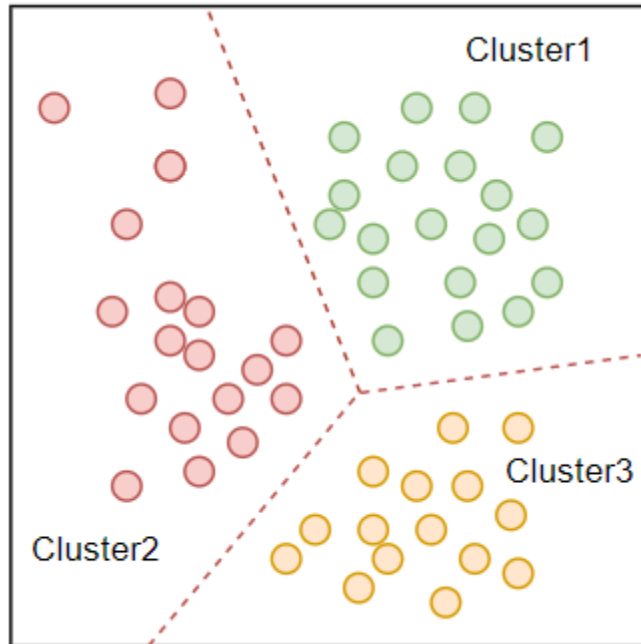




35 ^ 80

250 ^ 150

# Mastering the game of Go

**How to find the best local move for wining the whole game?**

(1) Deep neural networks supervised training

(2) Monte Carlo tree search programs

(3) Alpha-go Zero (2017) : only reinforcement from scratch



hidden layer 1   hidden layer 2   hidden layer 3

input layer

output layer

$p_{\sigma/\rho}\,(a|s)$

s

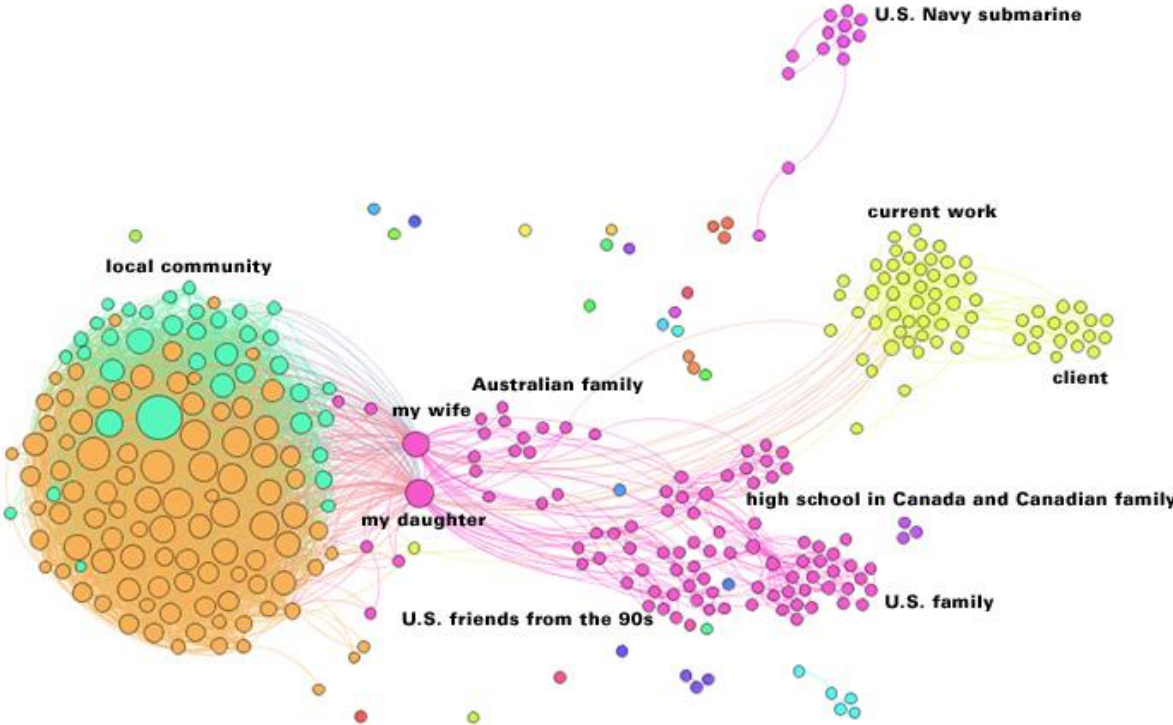# 2. Unsupervised learning (clustering)

Unsupervised ML learns from a dataset without any labels. The algorithm can automatically classify or categorize the input data.
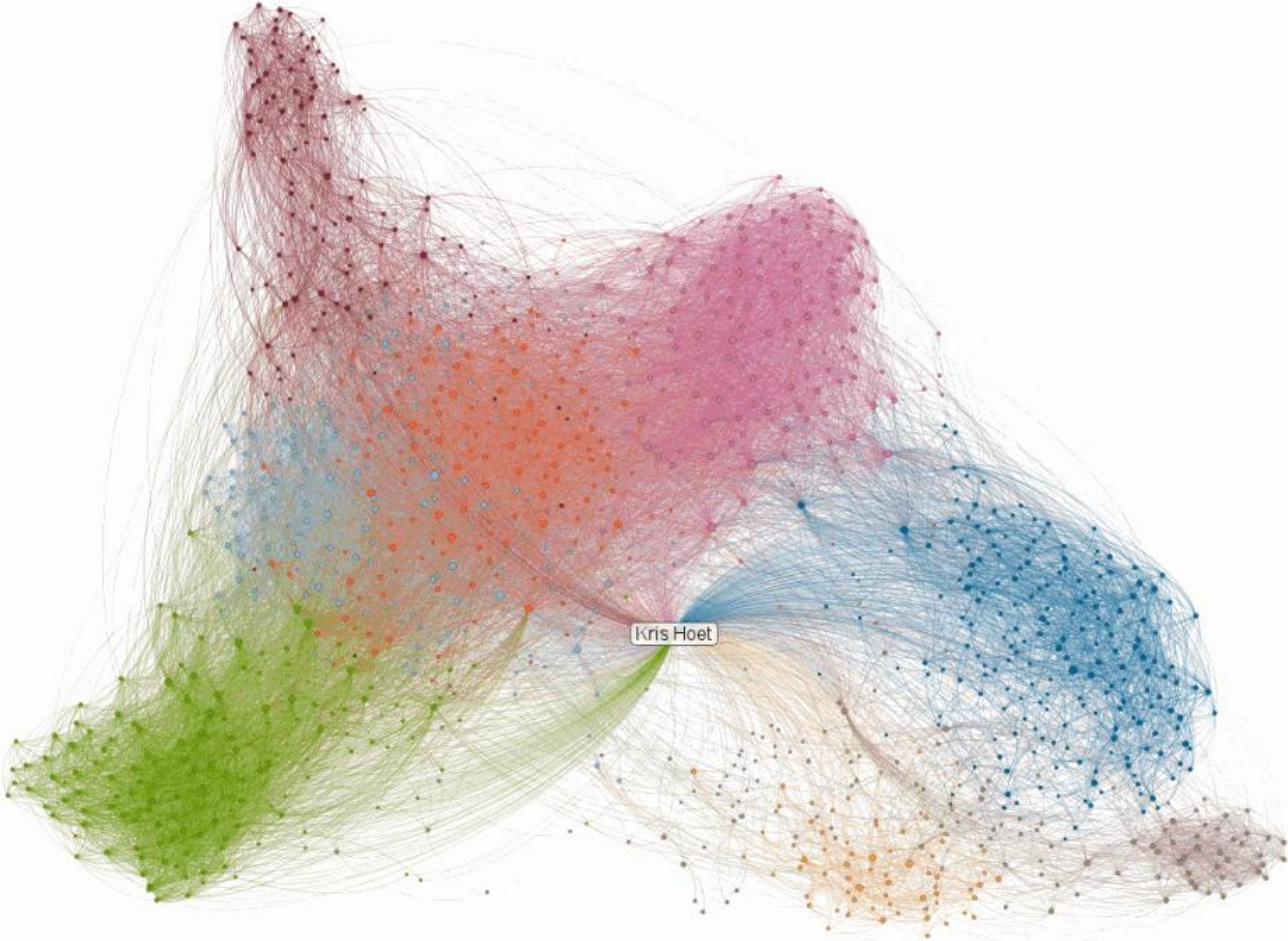


The application of unsupervised learning mainly includes cluster analysis, association rule or dimensionality reduce.

# Social Network Analysis
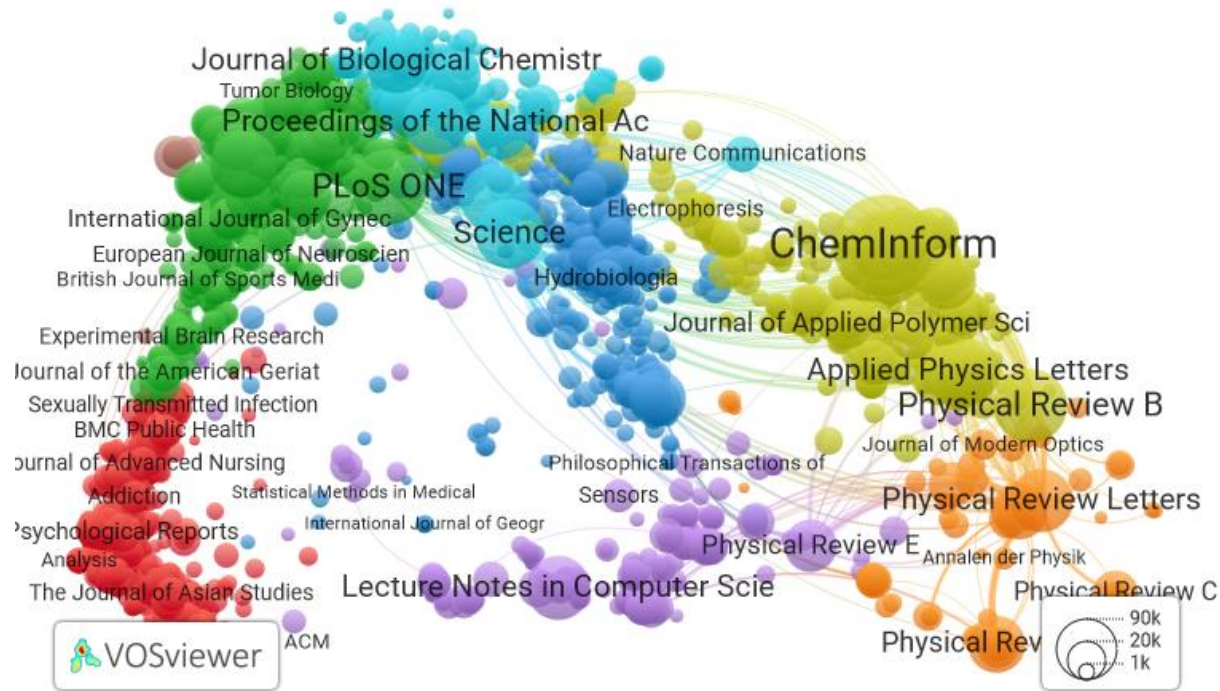


Chad's Facebook network October 2012



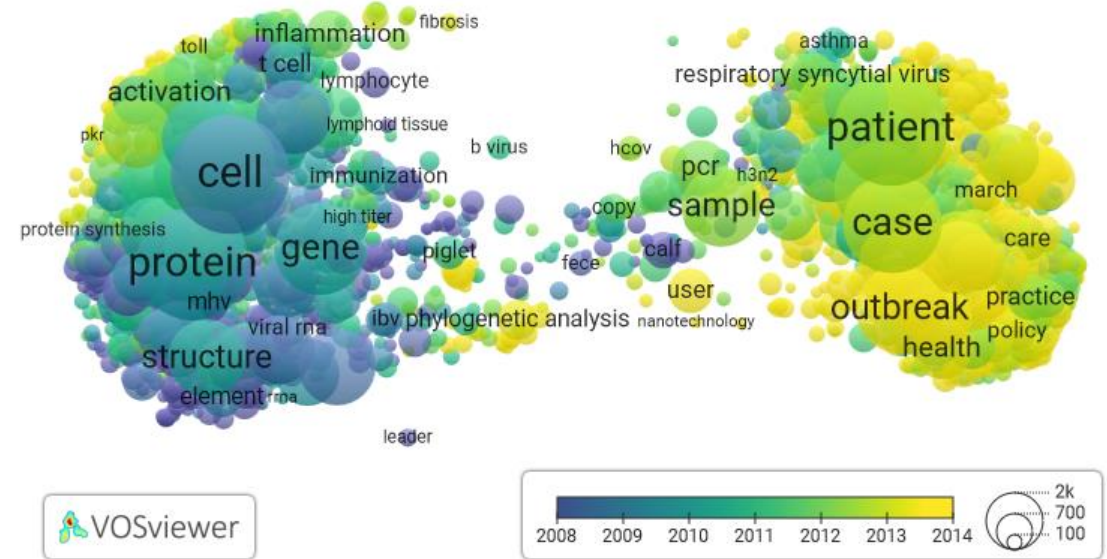LinkedIn. Maps — Kris Hoet's Professional Network as of March 1, 2011
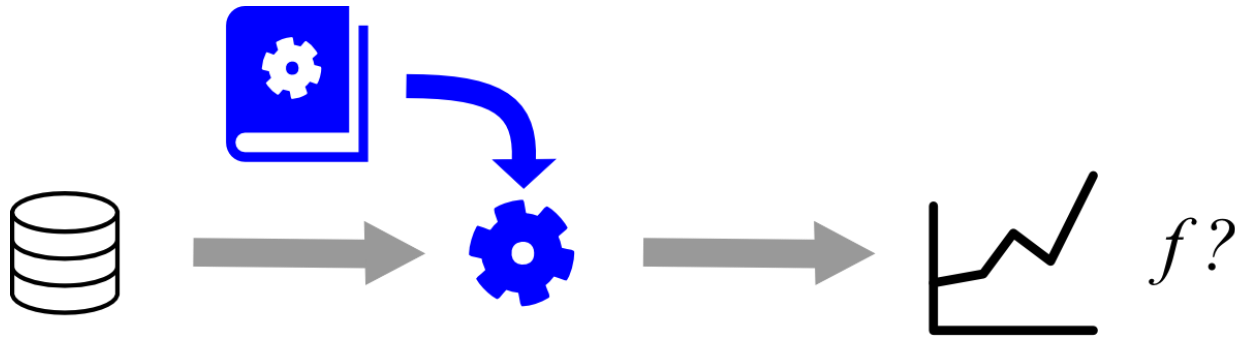
# Social Network Analysis



Citation network of journals



Co-occurrence network of terms in COVID-19 articles
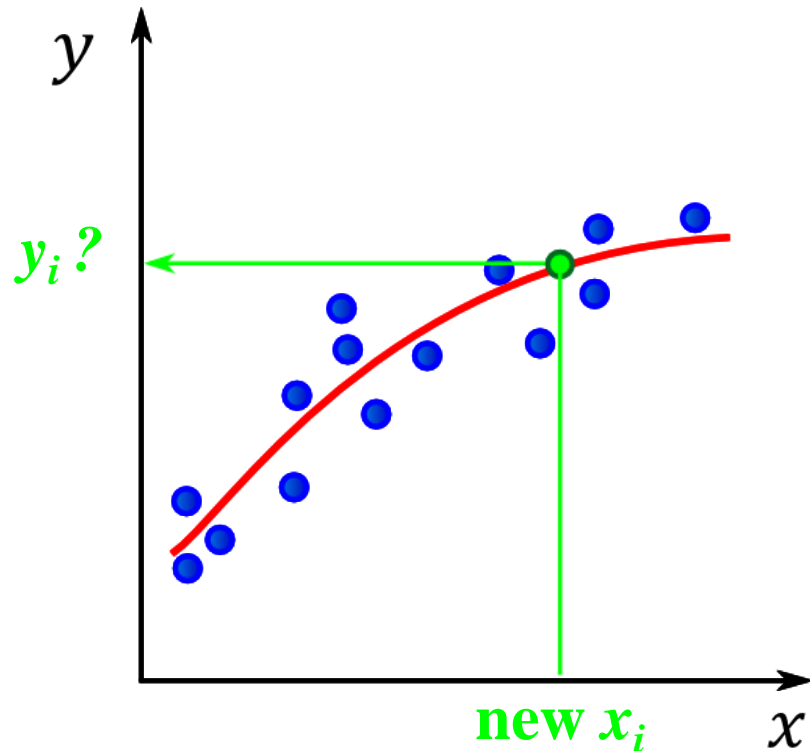
10

# 3. Supervised learning

Supervised ML learns from a trained tagged dataset, builds a function, predicts the output based on the function.
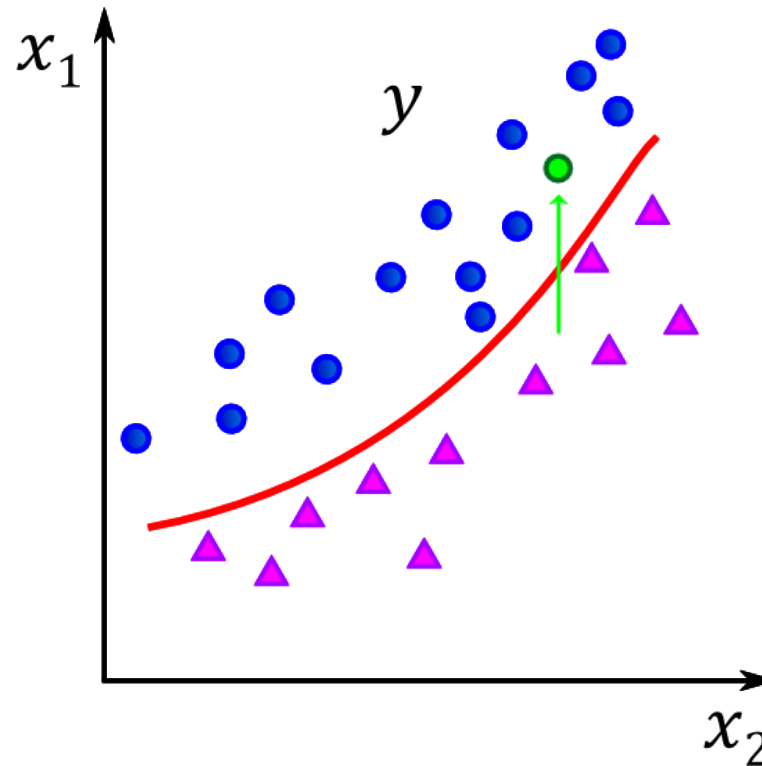
$f$? 

$f(\text{x})$ ?

(i) **quantitative variables:**
→ **Regression**

(ii) **qualitative variables:**
→ **Classification**

# 3. Supervised learning

The training dataset often consists of pairs of an input vector (or scalar) and the corresponding output vector (or scalar), the output of the function can be **regression** or **classification**.
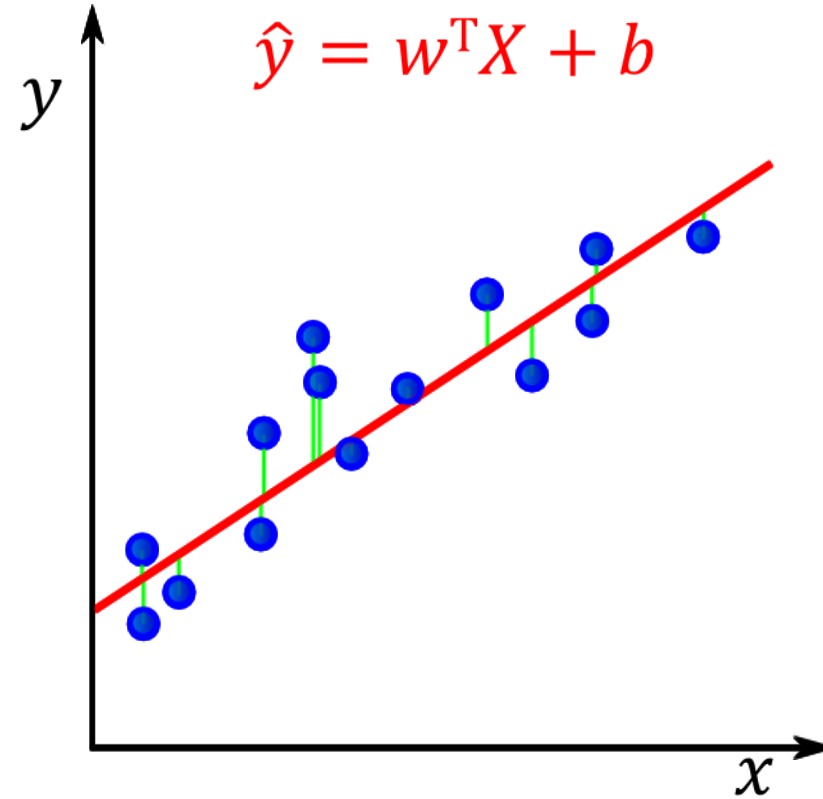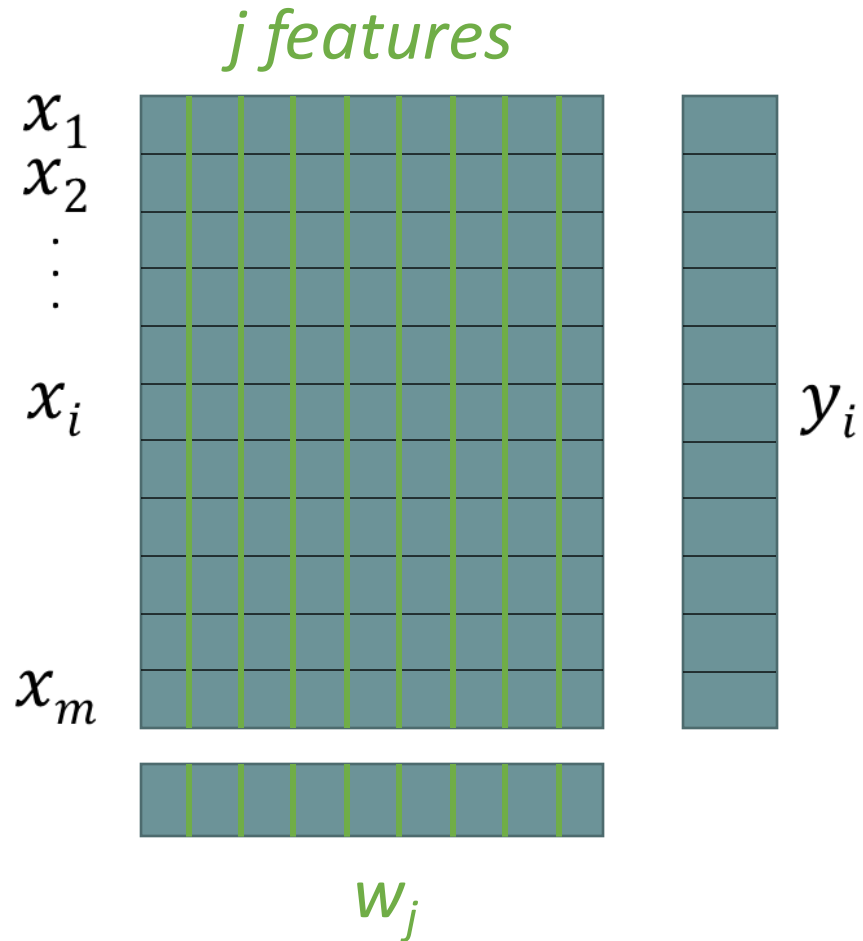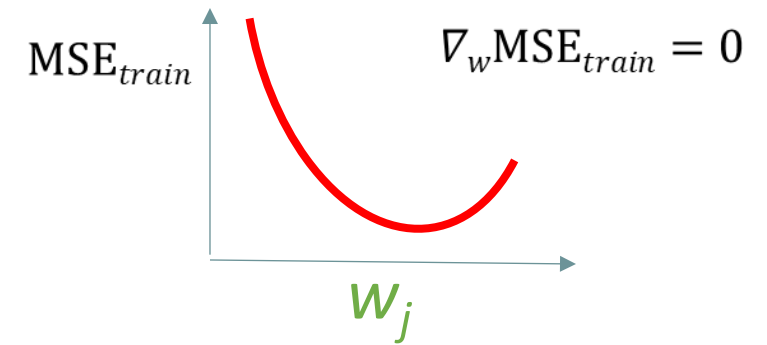


**Regression**  **Classification**

# Optimization of the learning

*j features*

$x_1$
$x_2$
⋮
$x_i$

$x_m$

$w_j$

$y_i$

$$\hat{y} = w^T X + b$$

$y$

$x$

$$SSR = \sum_{i}^{m} (\hat{y}_i - yi)^2$$

$$R^2 = 1 - \sum \frac{(\hat{y}_i - yi)^2}{(\bar{y}_i - yi)^2}$$

# Estimation of the learning



$x_1$
$x_2$
$\vdots$
$x_i$
$x_m$

$y$

training set

testing set

$w_j$

$$\mathrm{MSE}_{train}$$

$$\nabla_w \mathrm{MSE}_{train} = 0$$

$w_j$

$$\mathrm{MSE}_{test} = \frac{1}{m}\sum_{i}^{m}(\widehat{y_{test}} - yte_{st})^2$$

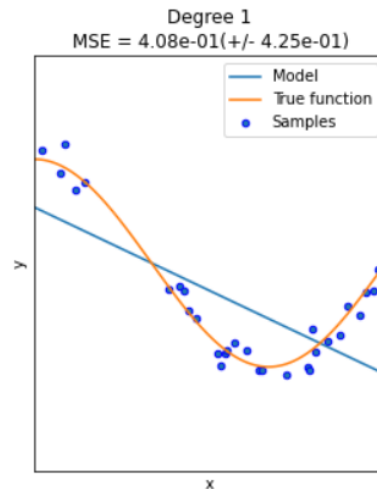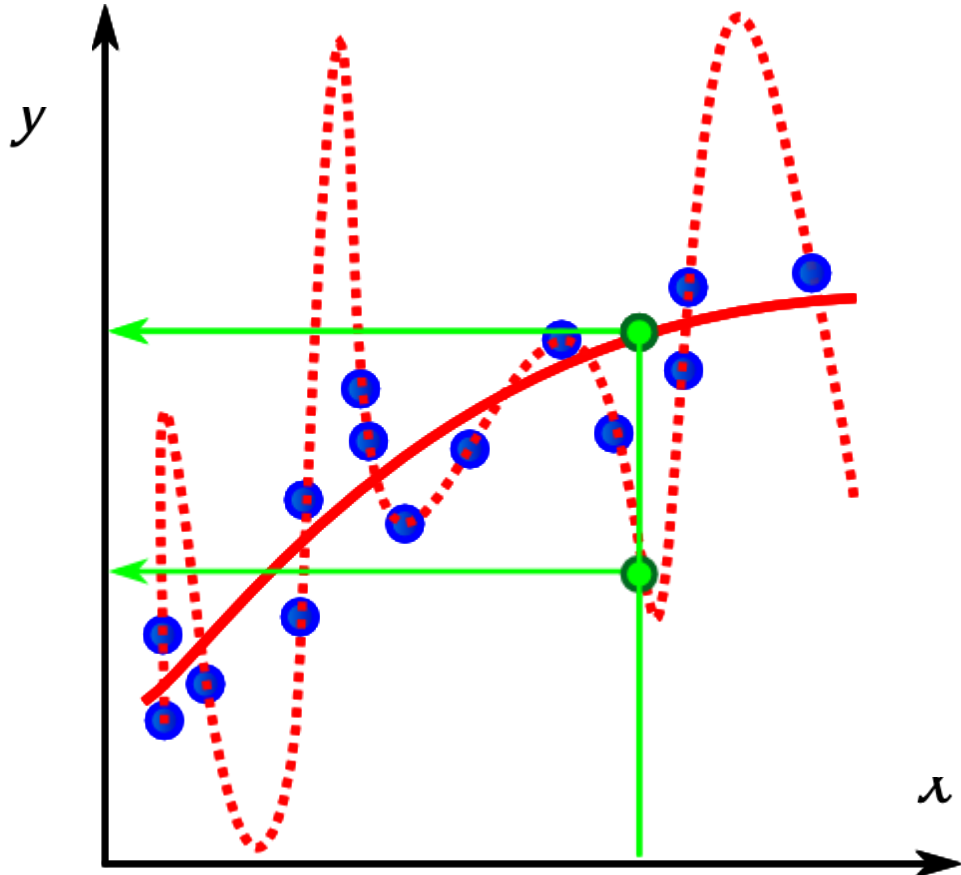$$\mathrm{RMSE}_{test} = \sqrt{\mathrm{MSE}_{test}}$$

# Cross Validation



If large dataset, CV is need by repeating training/testing procedure under *k*-folder: partitions formed by splitting into *k* non-overlapping subsets.
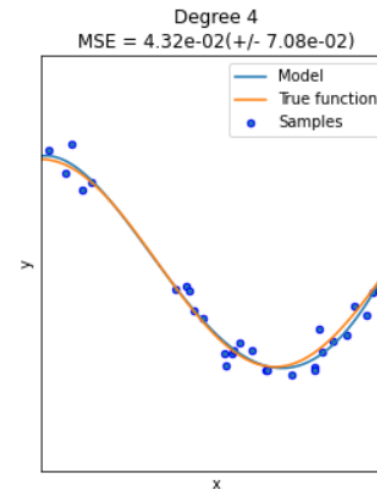
# Overfitting

= the production of an analysis that corresponds too closely or exactly to a particular set of data, and may therefore fail to fit additional data or predict future observations reliably.
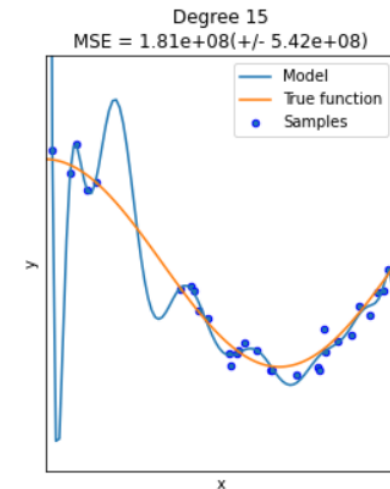
$$\hat{y} = \sum_{i}^{n} w_i X^i + b$$

Degree 1
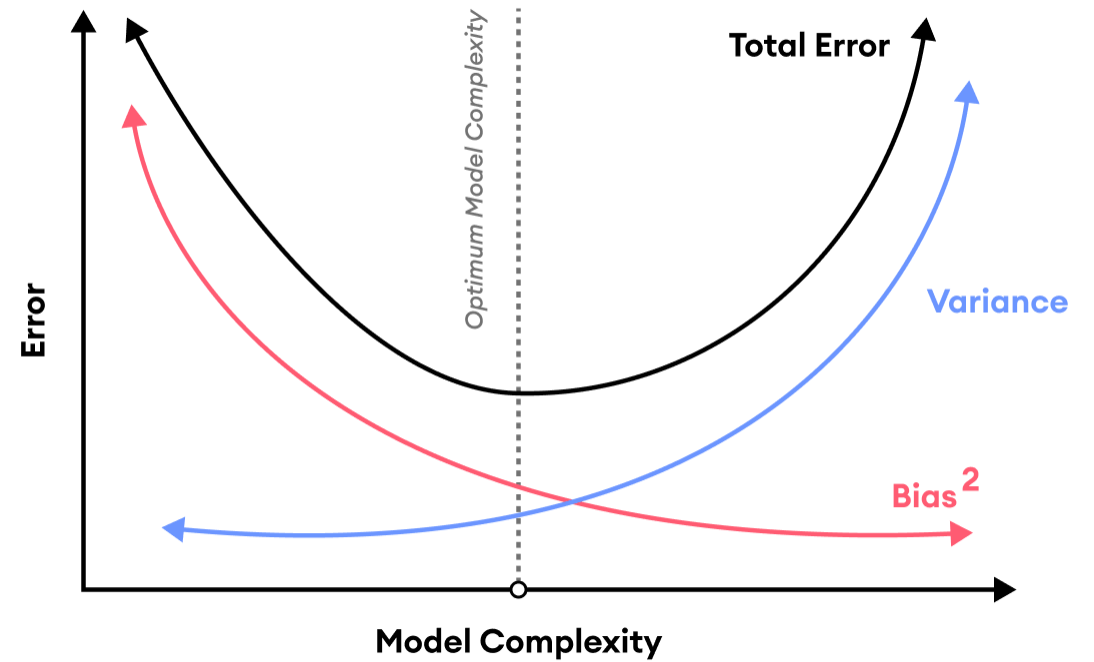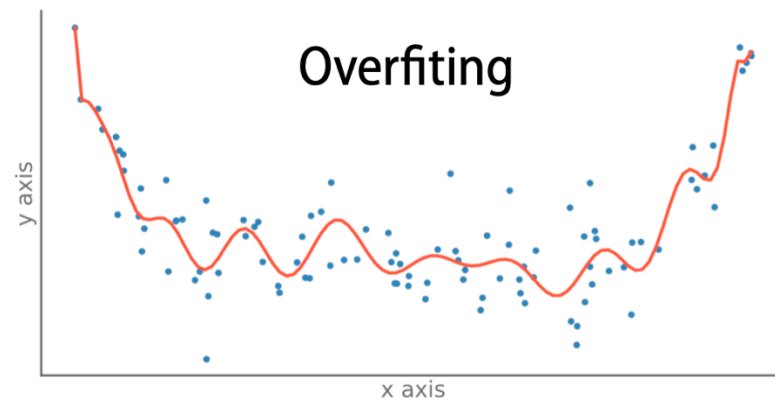MSE = 4.08e-01(+/- 4.25e-01)

- Model
- True function
- Samples

Degree 4
MSE = 4.32e-02(+/- 7.08e-02)

- Model
- True function
- Samples

Degree 15
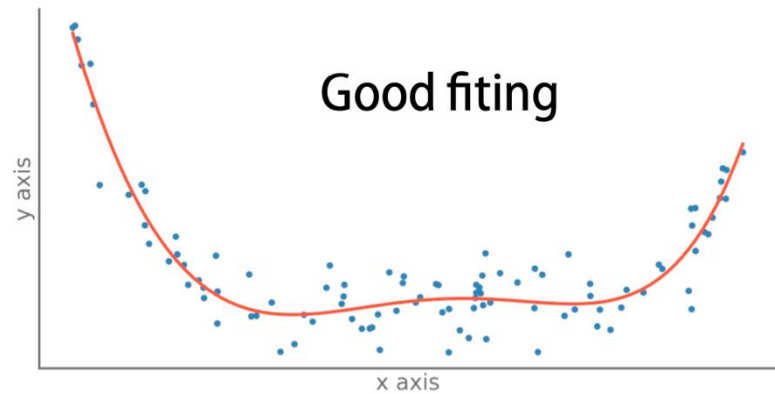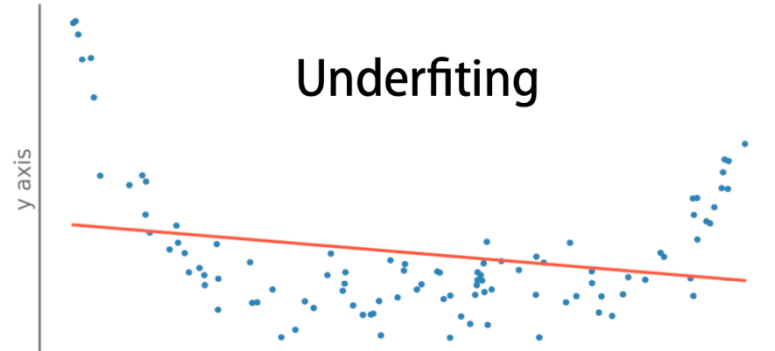MSE = 1.81e+08(+/- 5.42e+08)

- Model
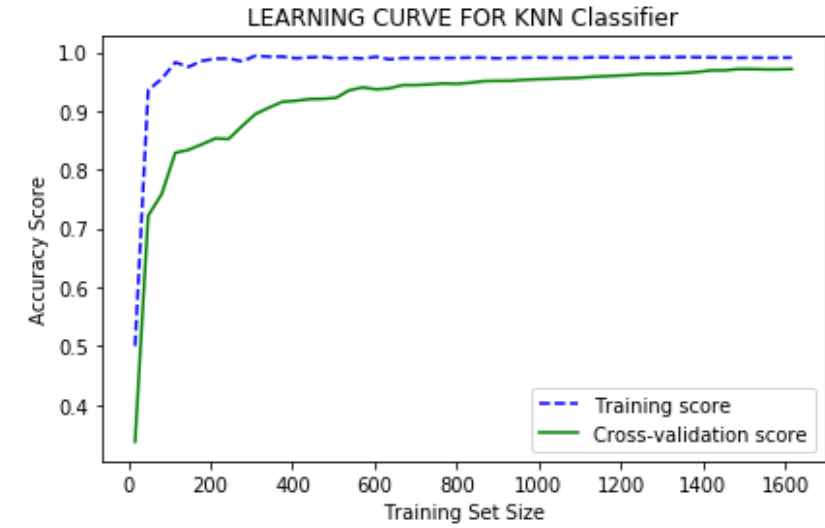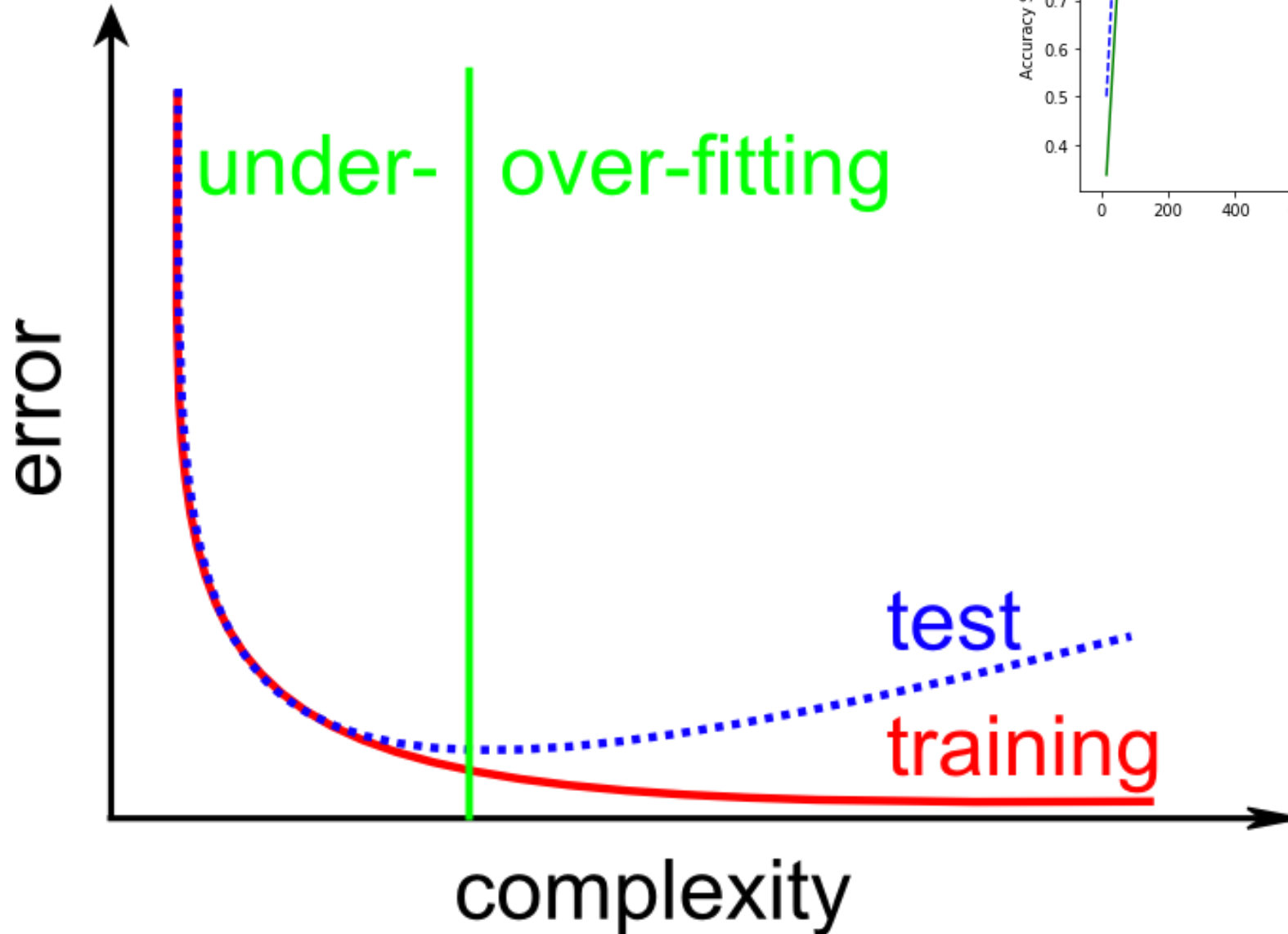- True function
- Samples

high biais
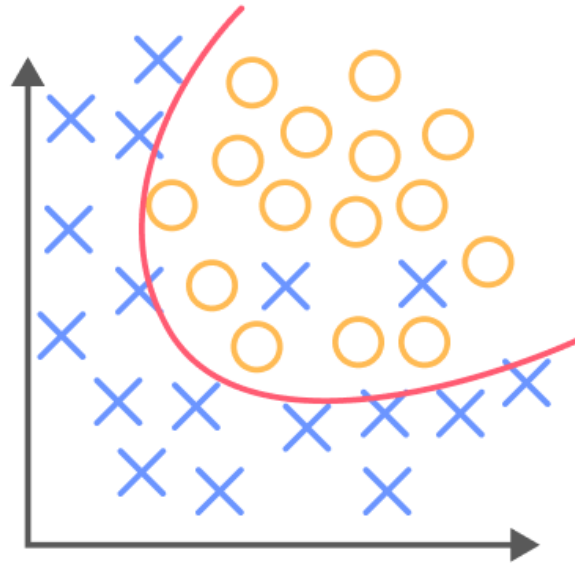underfitting

just right

low biais
overfitting

# Bias-variance tradeoff

# Overfitting



under- | over-fitting

error

complexity

test

training

LEARNING CURVE FOR KNN Classifier
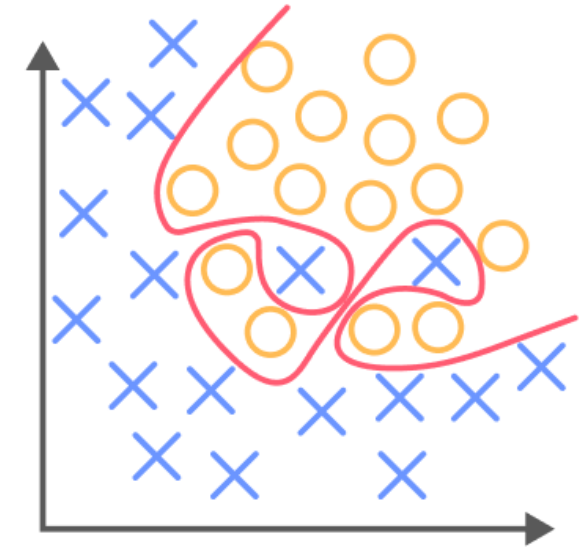
Accuracy Score

Training Set Size

- - - Training score
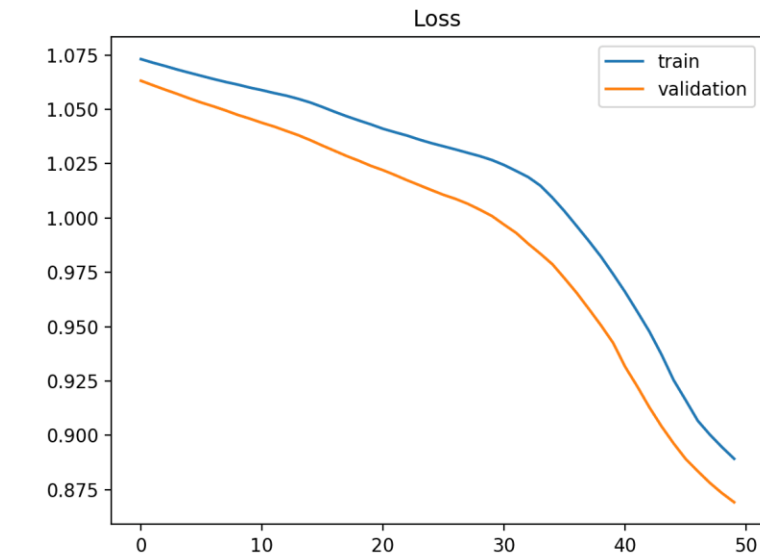— Cross-validation score

# Overfitting in classification

Underfitting
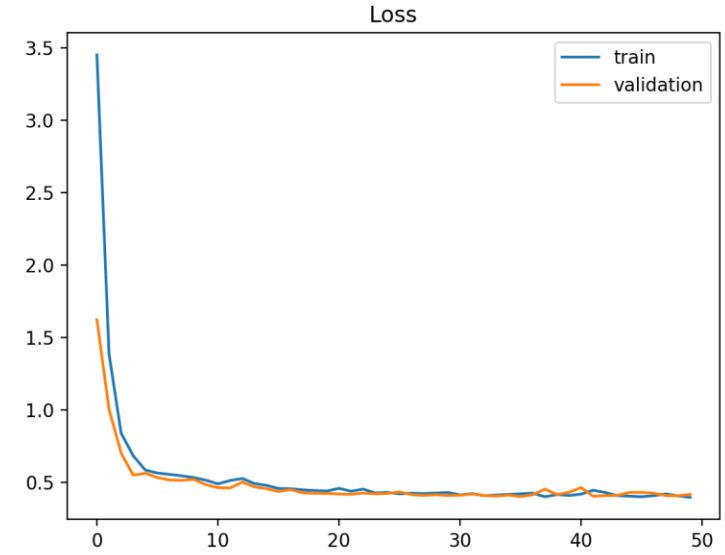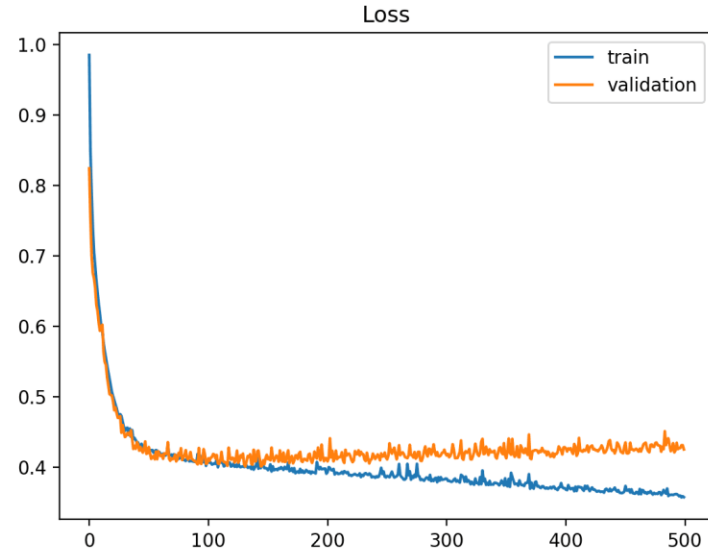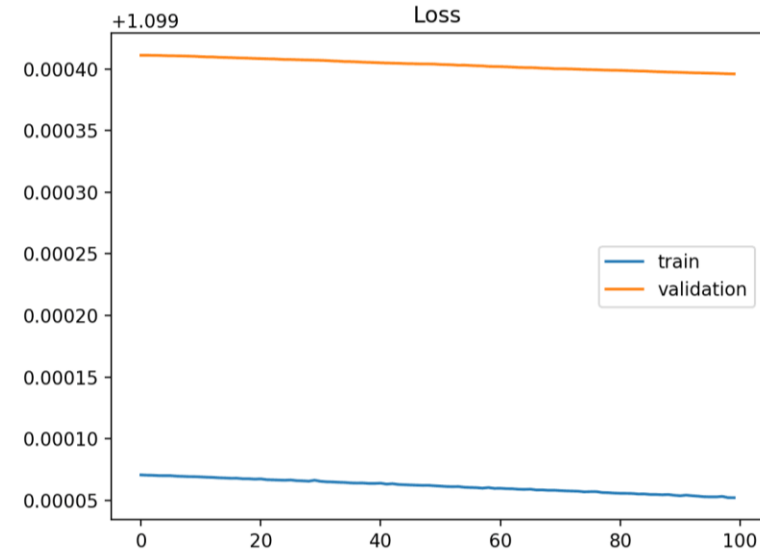
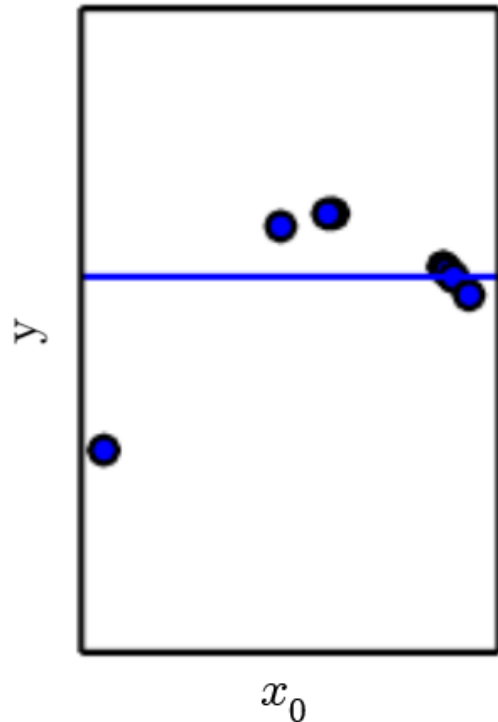Appropriate fitting

Overfitting

# Learning curve: performance VS time
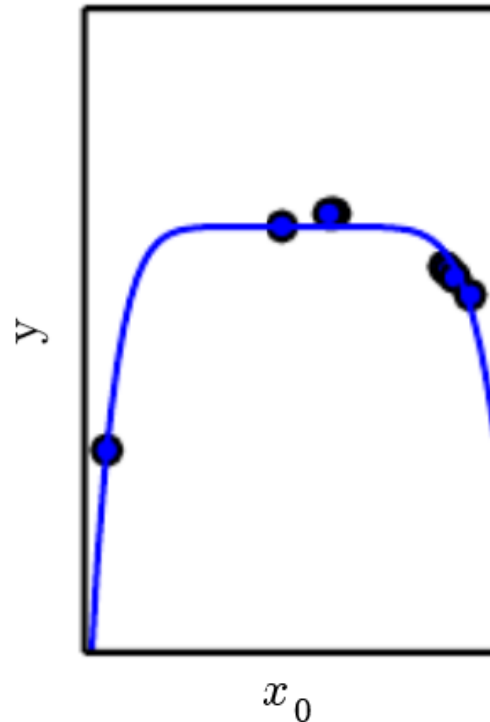
# Regularization by penalties

We can introduce a weight decay to degrade the learning.

$$\text{error}(w) = \text{MSE}_{train} + \lambda w^{\text{T}} w$$
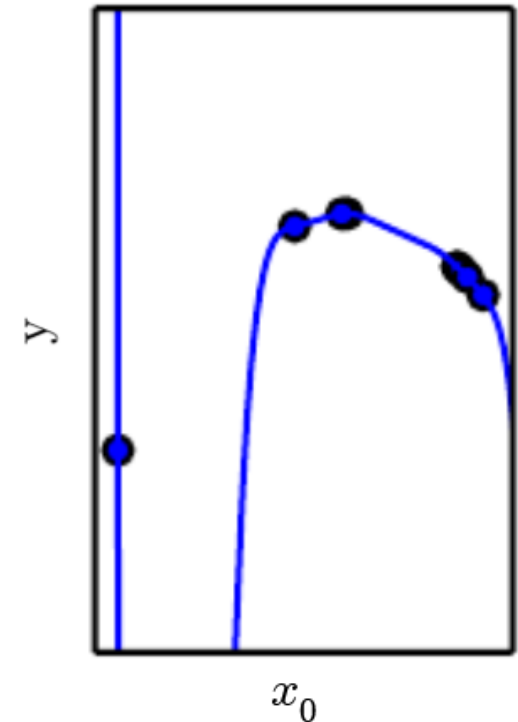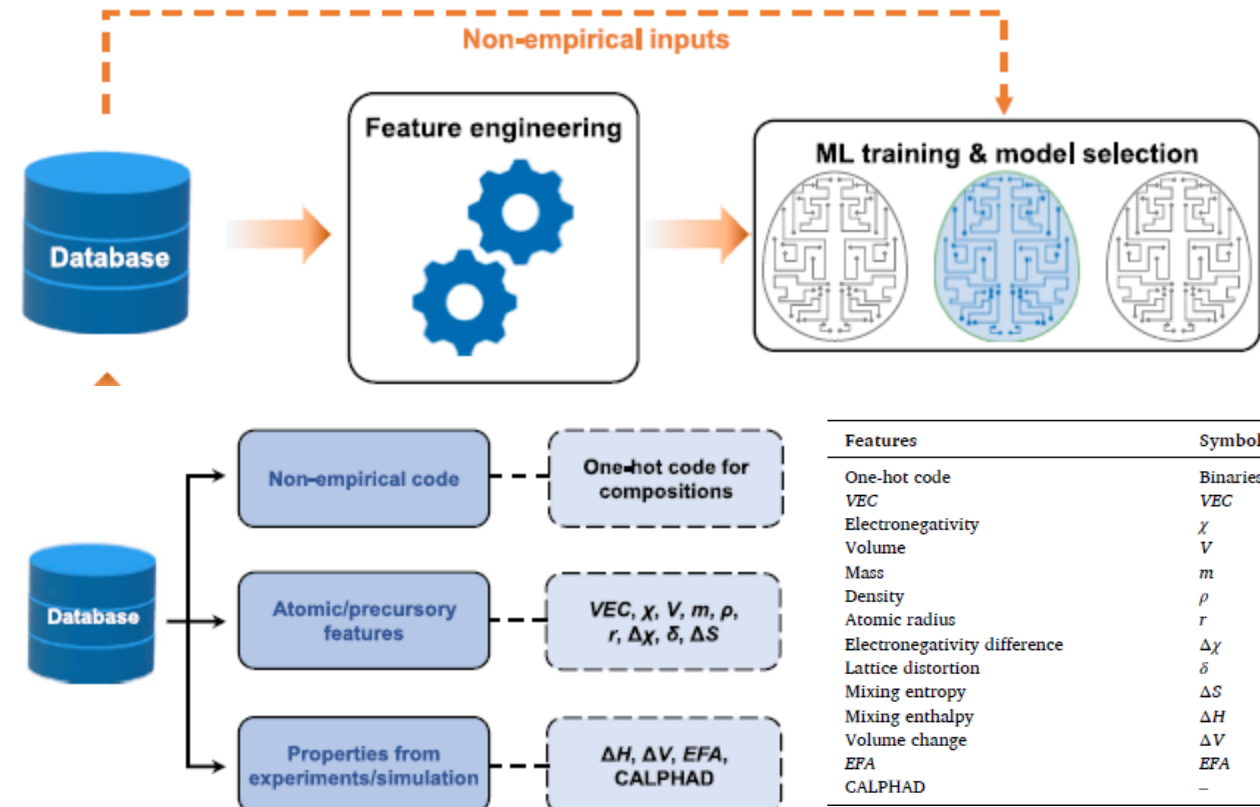


Underfitting
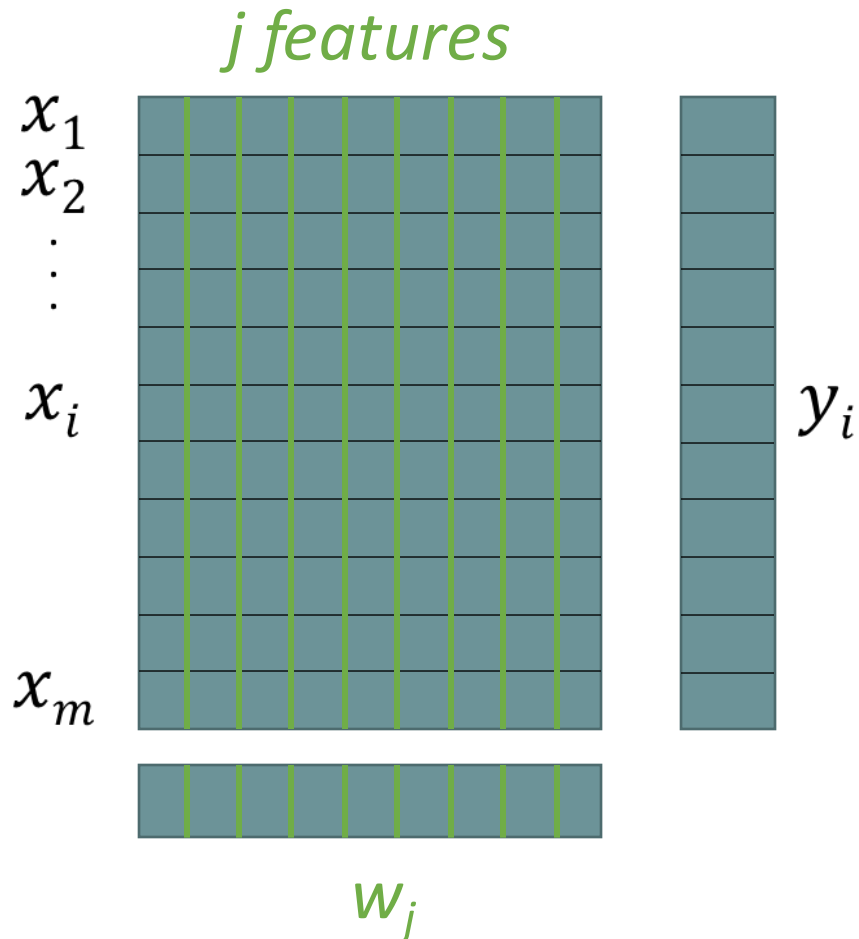(Excessive $\lambda$)

Appropriate weight decay
(Medium $\lambda$)

Overfitting
($\lambda \rightarrow 0$)

# Choice of descriptors / features



*j features*

$x_1$
$x_2$
$\vdots$
$x_i$
$x_m$

$y_i$

$w_j$



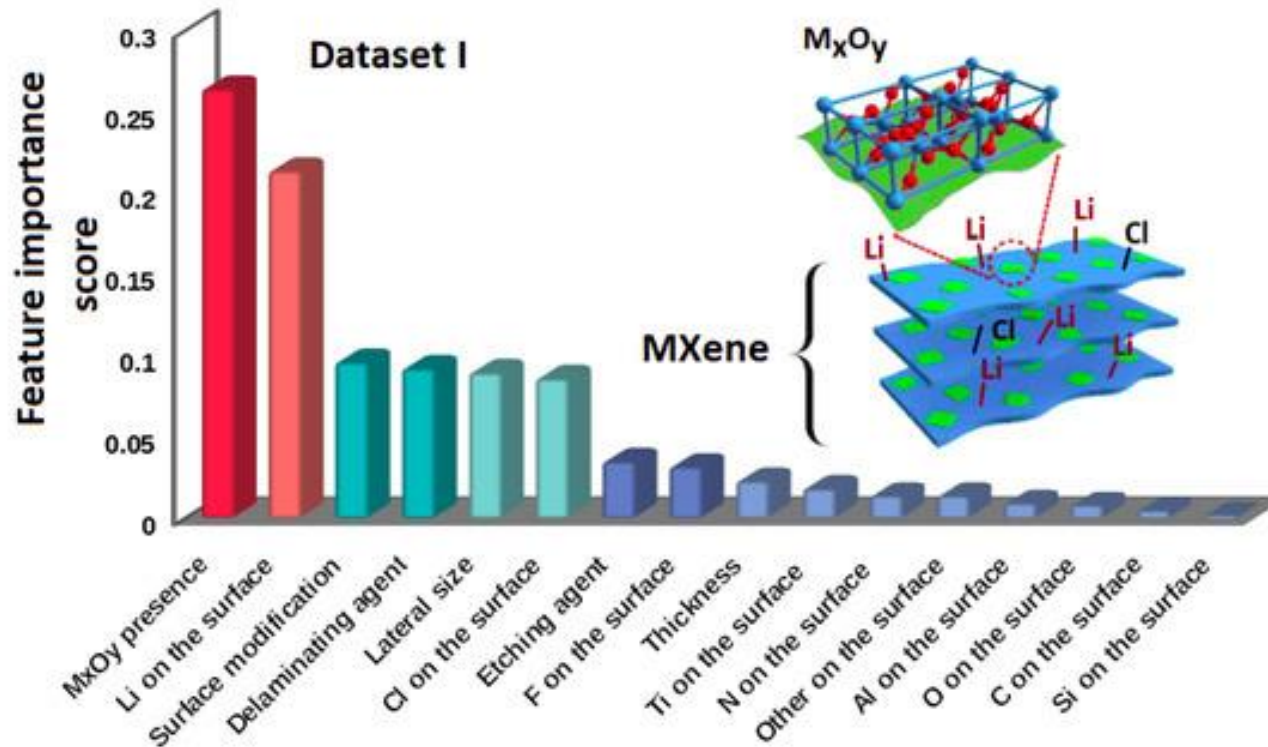**Zhang *et al*. Curr. Opin. Solid State Mater. Sci. (2020)**
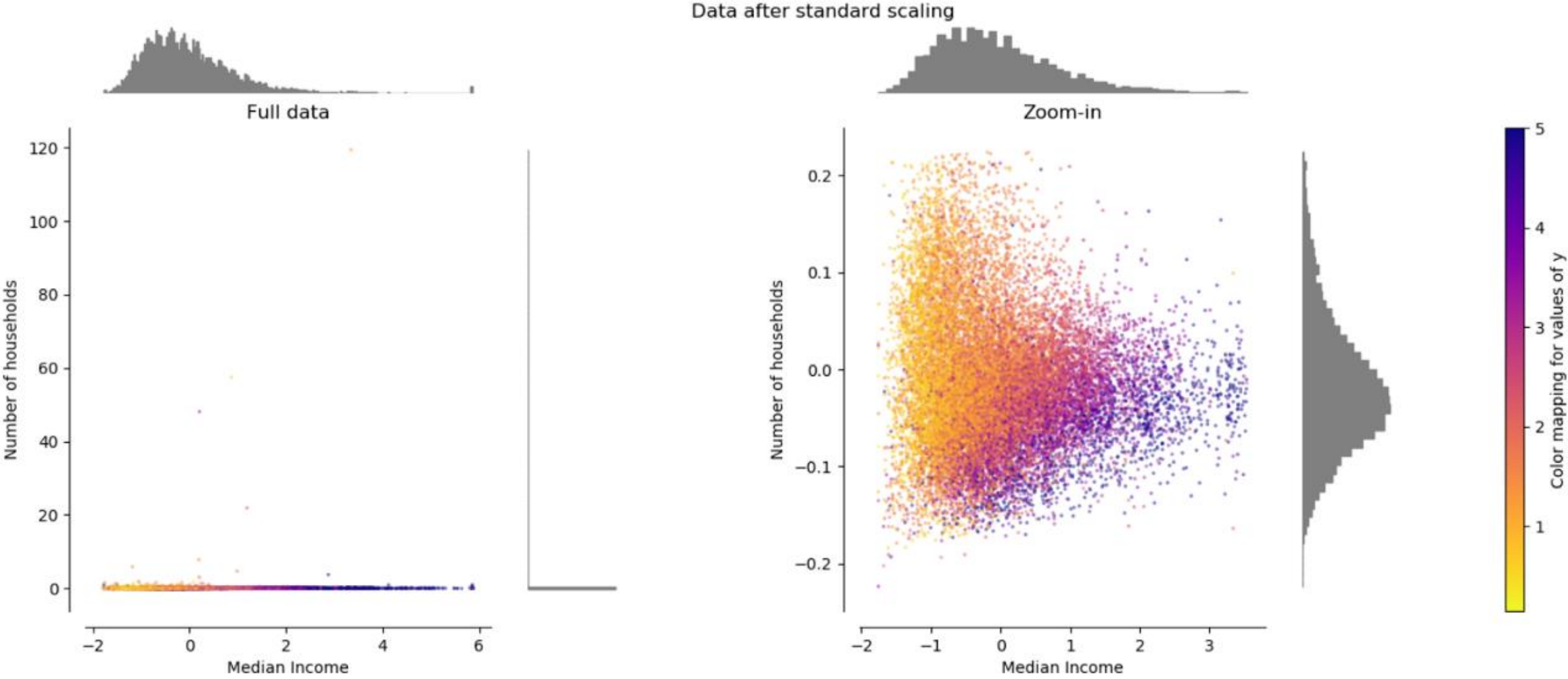*Rational design of high-entropy ceramics based on machine learning – A critical review*

# Feature importance



**Marchwiany *et al*.   Materials (2020)**
*Surface-Related Features Responsible for Cytotoxic Behavior of MXenes Layered Materials Predicted with Machine Learning Approach*

# Standardization

# Standardization



Actual Data · After normalizing · After standardization

Moyenne: $\mu = \frac{1}{N}\sum_{i=1}^{N} x_i$

Variance: $V = \frac{1}{N}\sum_{i=1}^{N}(x_i - \mu)^2$

Ecart type: $\sigma = \sqrt{V}$

$$z_i = \frac{x_i - \mu}{\sigma}$$

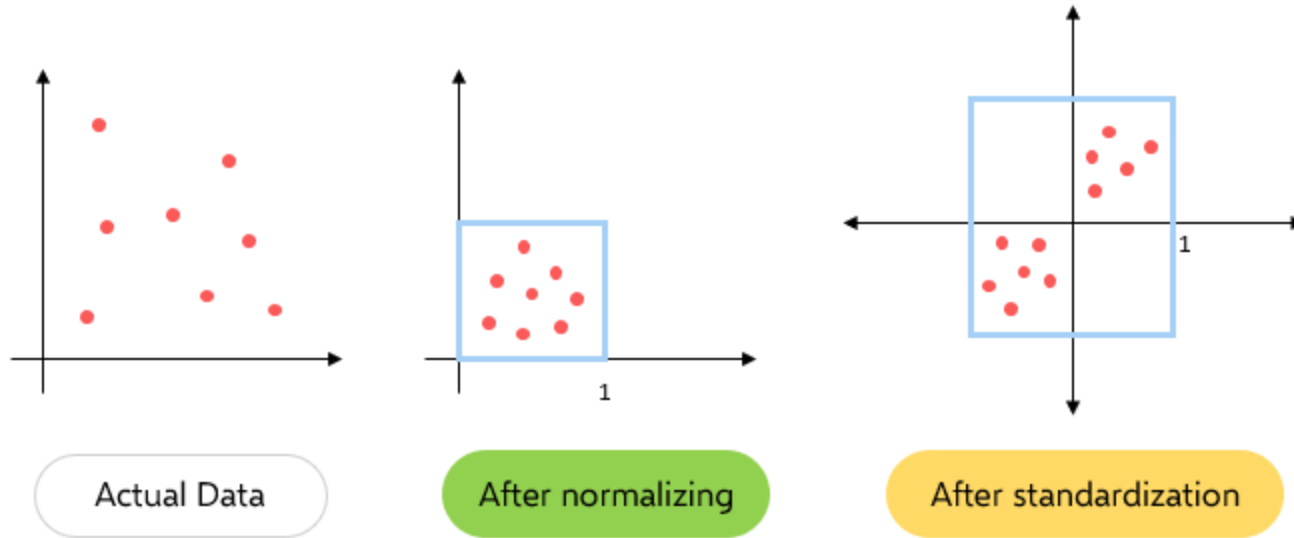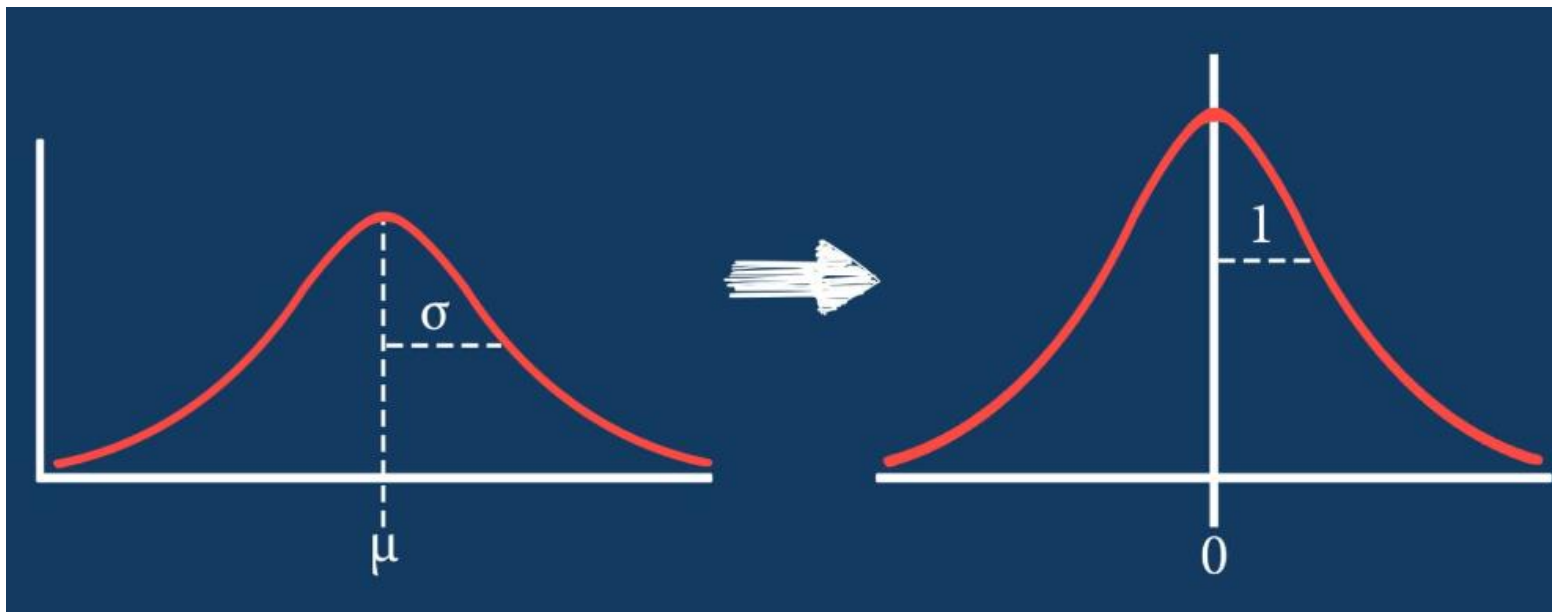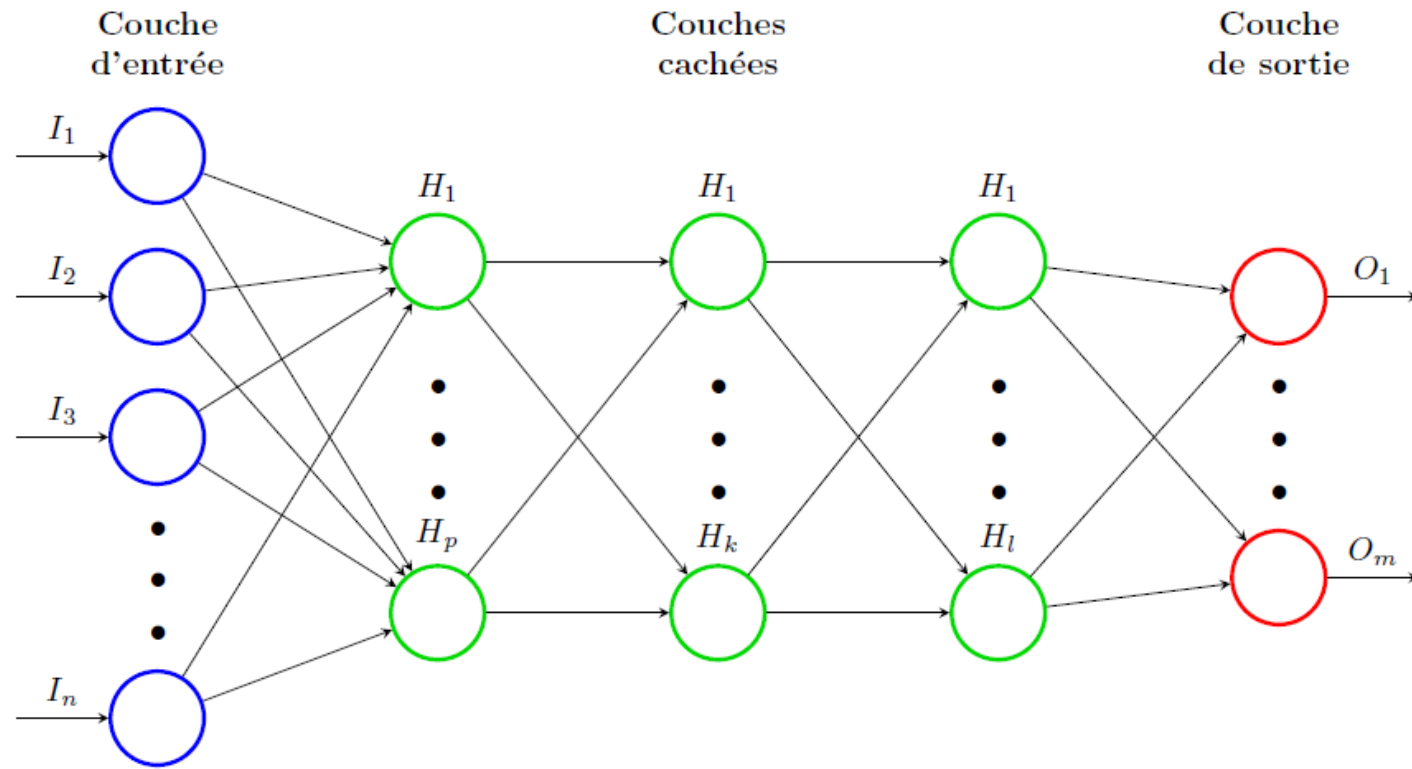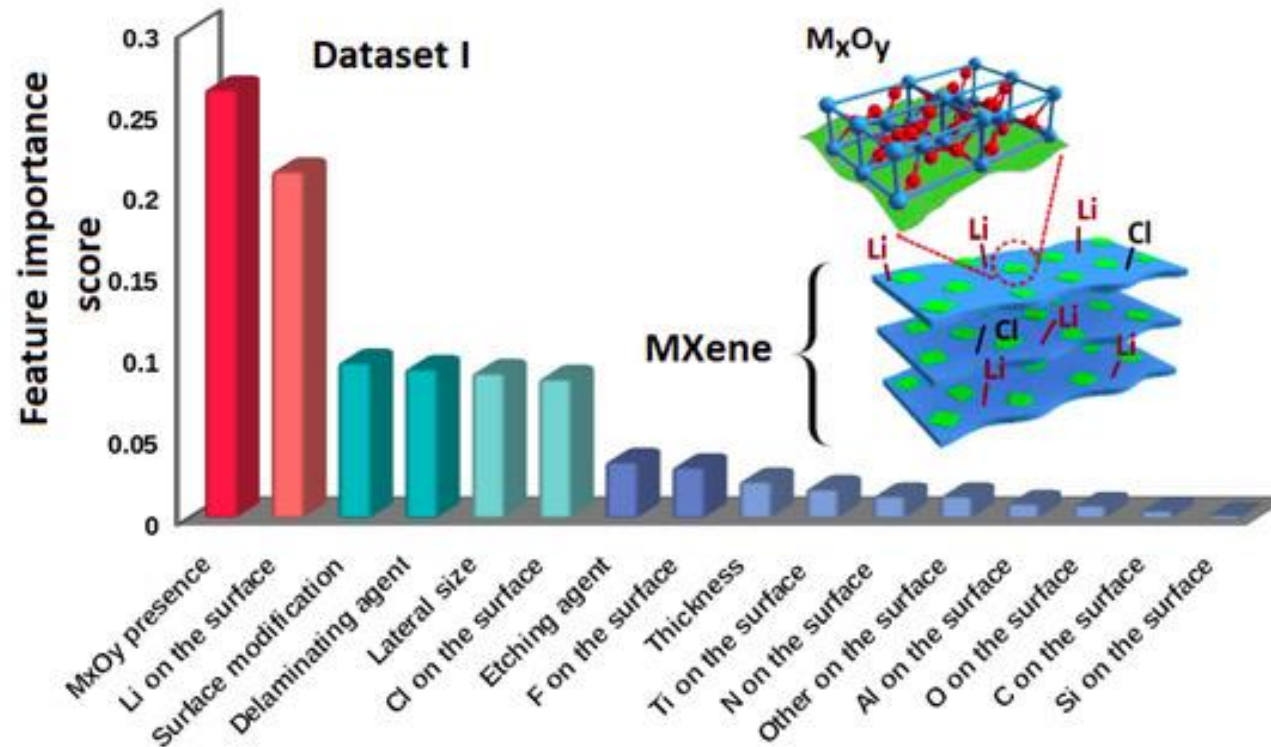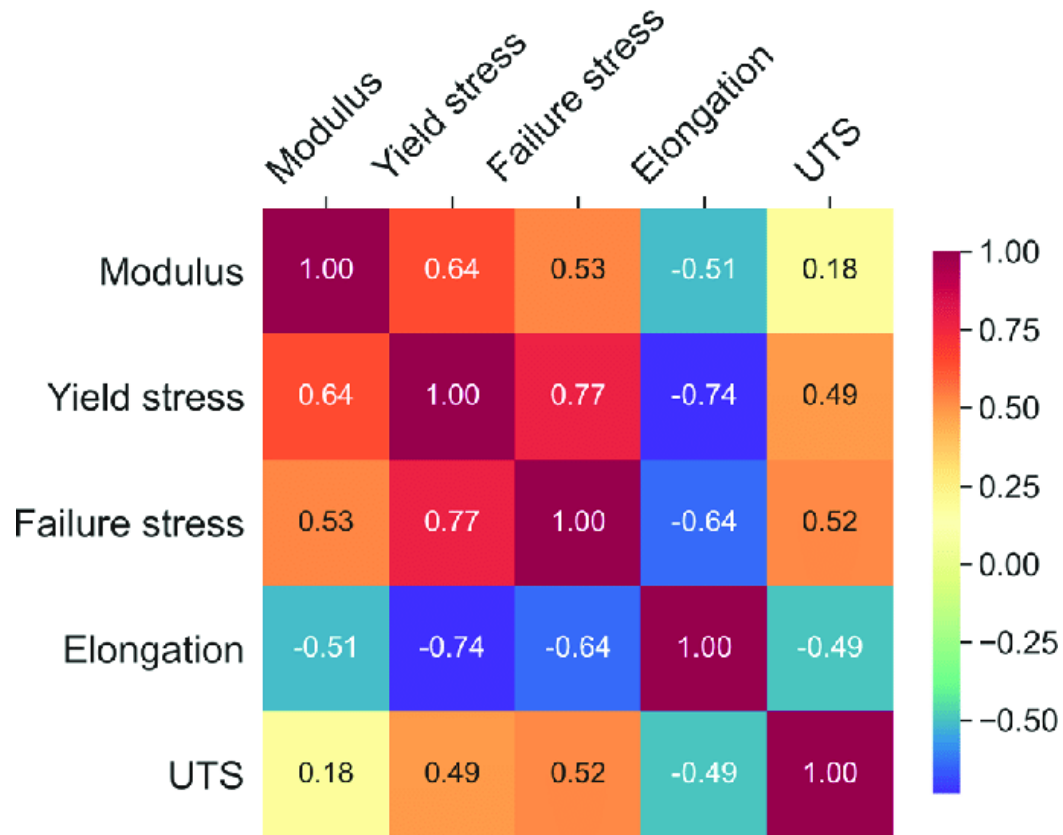# Why is standardisation important?

# Feature importance



**Marchwiany *et al*.  Materials (2020)**
*Surface-Related Features Responsible for Cytotoxic Behavior of MXenes Layered Materials Predicted with Machine Learning Approach*

# Pair correlation matrix



Esperance: $E[X] = \sum_{i=1}^{\infty} x_i\, p_i$
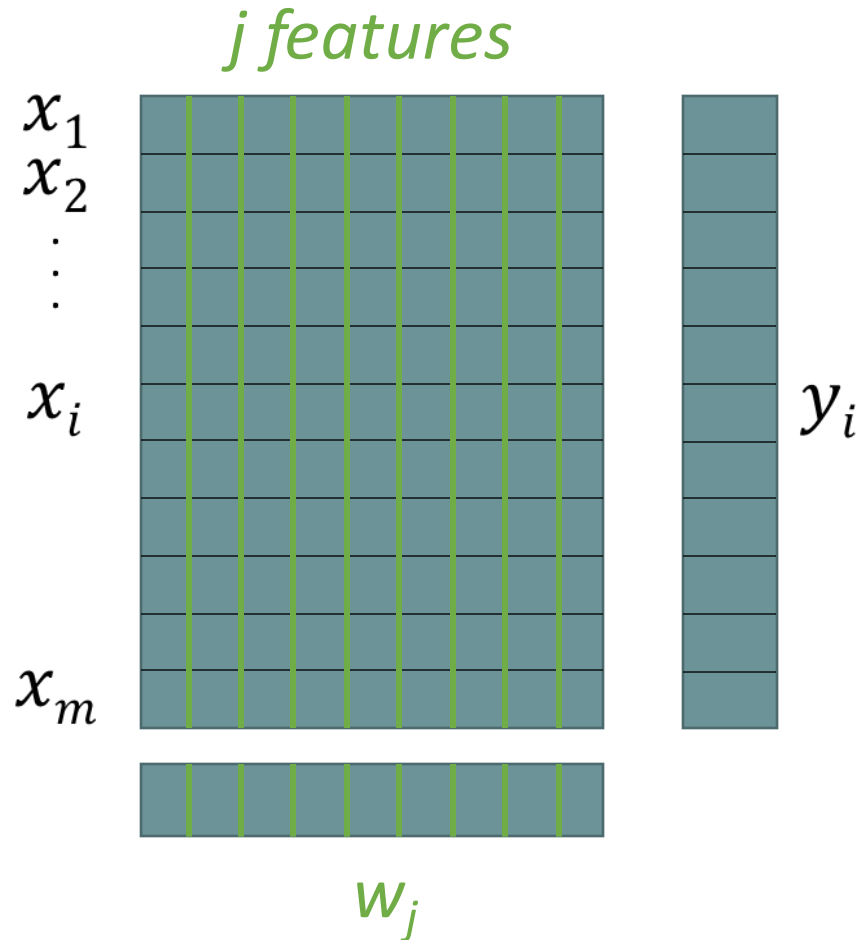
Covariance: $\mathrm{Cov}(X, Y) =$

$$\sum_i \sum_j x_i y_j\, \mathrm{P}(X = xi \text{ et } Y = yj) - E[X]E[Y]$$

# One-Hot Encoding the categorical variables

otherwise known as dummy variables, is a method of converting categorical variables into several binary columns
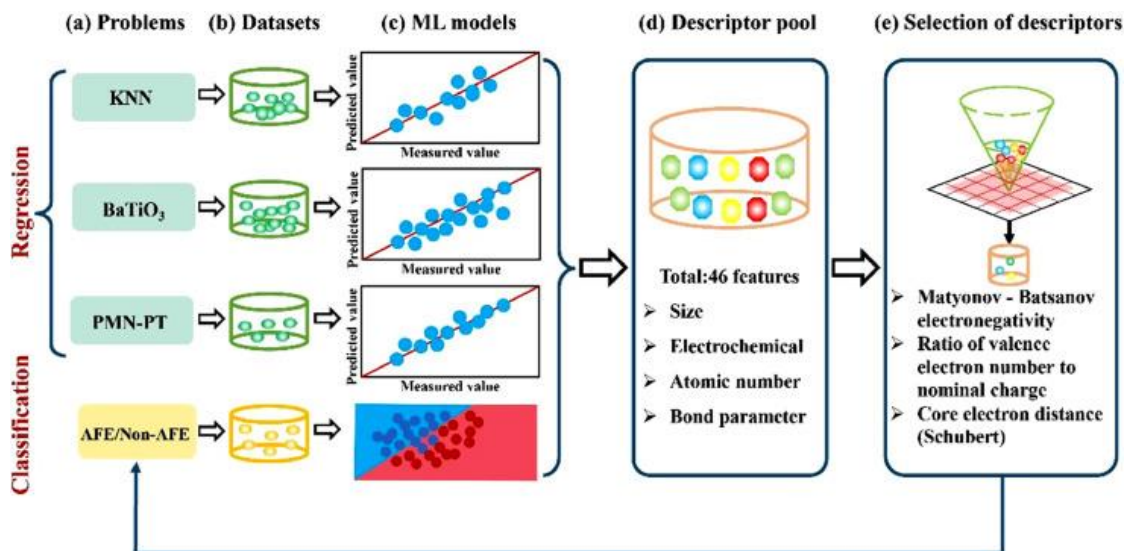
*j features*

$x_1$
$x_2$
$\vdots$
$x_i$

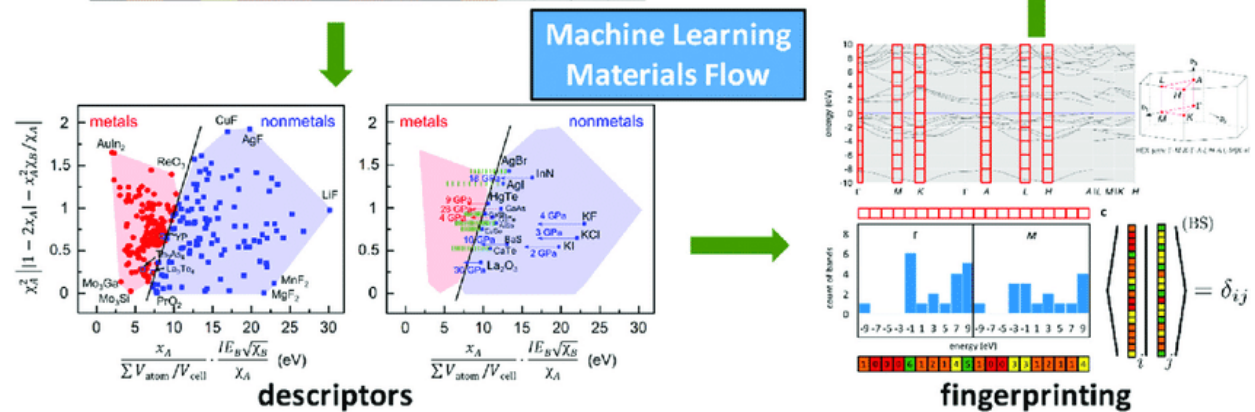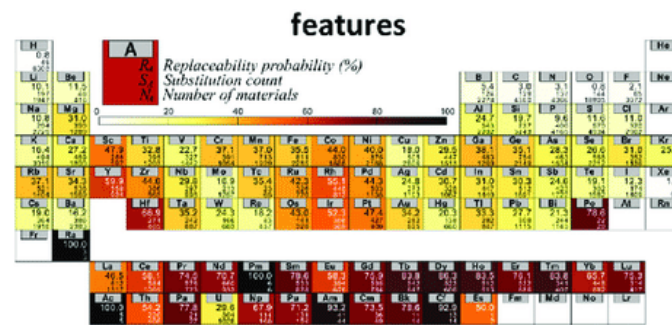$y_i$

$x_m$

$w_j$

Human-Readable

| Pet |
|-----|
| Cat |
| Dog |
| Turtle |
| Fish |
| Cat |

Machine-Readable

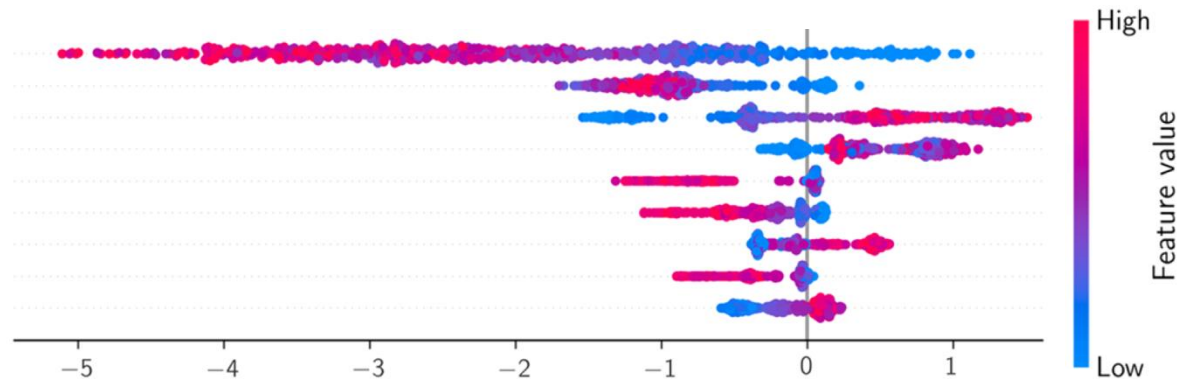| Cat | Dog | Turtle | Fish |
|-----|-----|--------|------|
| 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 |
| 1 | 0 | 0 | 0 |

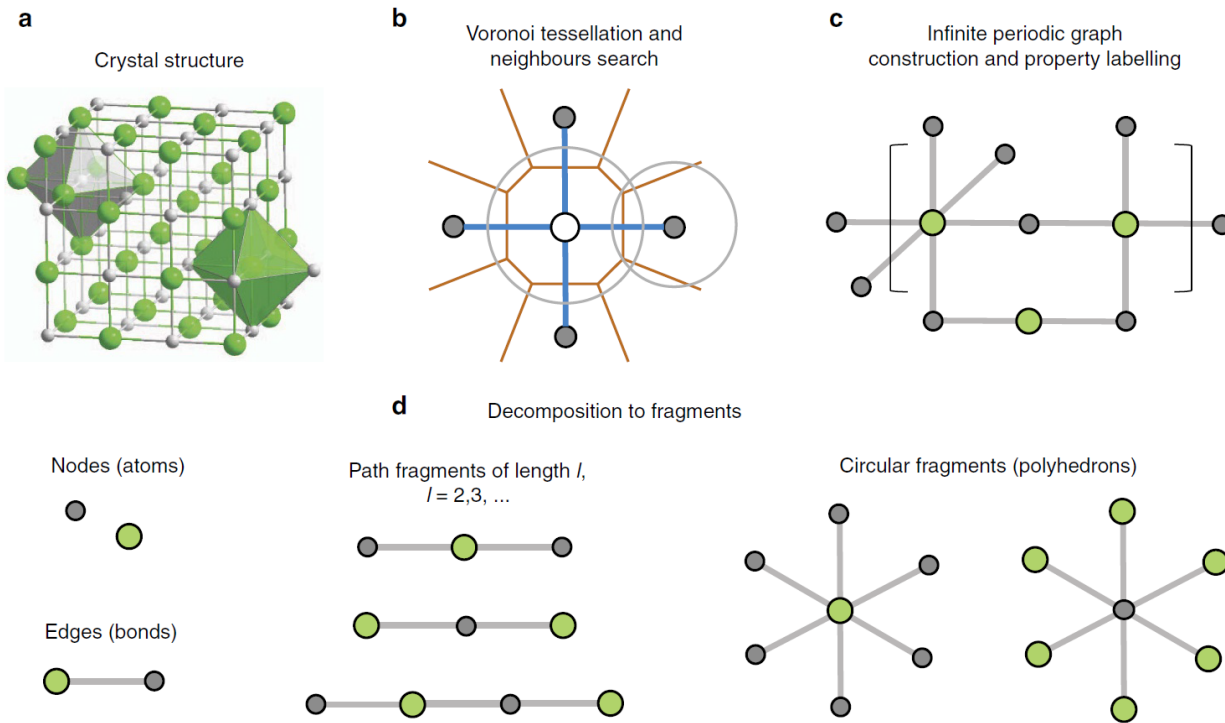# Choice of descriptors

# Improved performance with embedded physics



**Witman *et al.***
**Chem Mater (2021)**
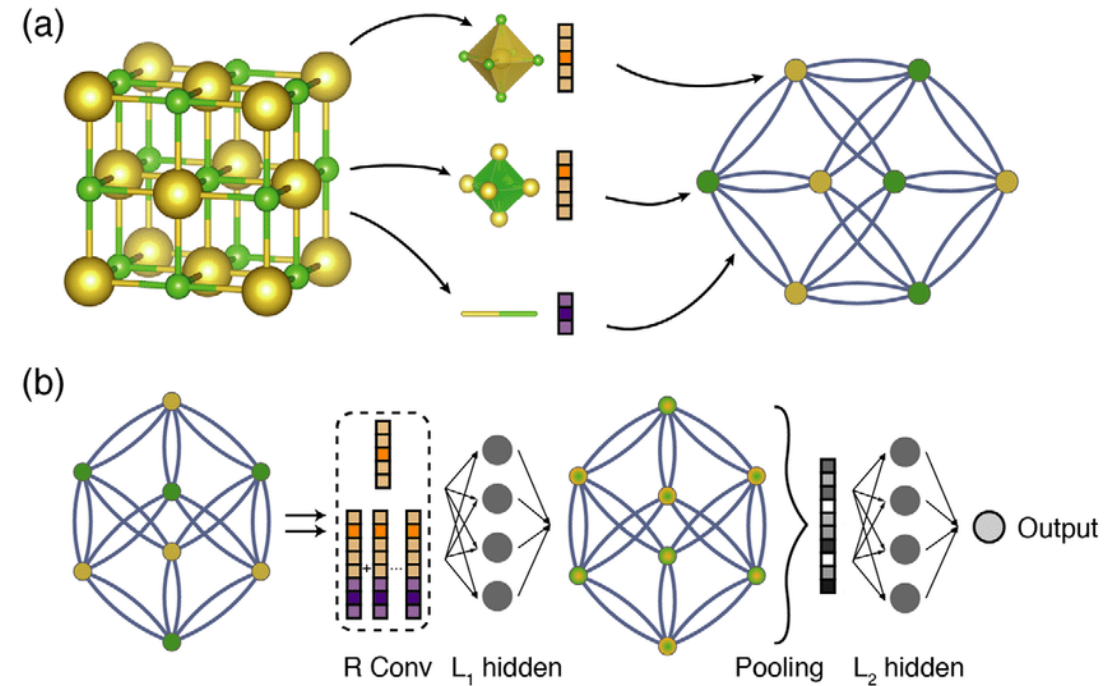*Data-Driven Discovery and Synthesis of High Entropy Alloy Hydrides with Targeted Thermodynamic Stability*
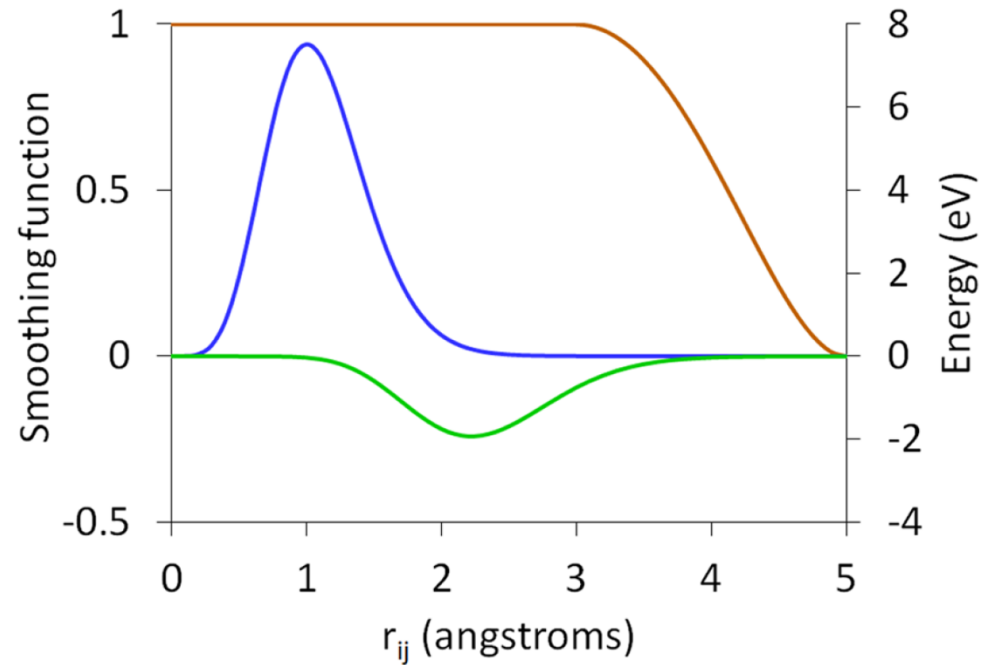
# Graphs description



Isayev *et al.*
**Nature Com (2017)**
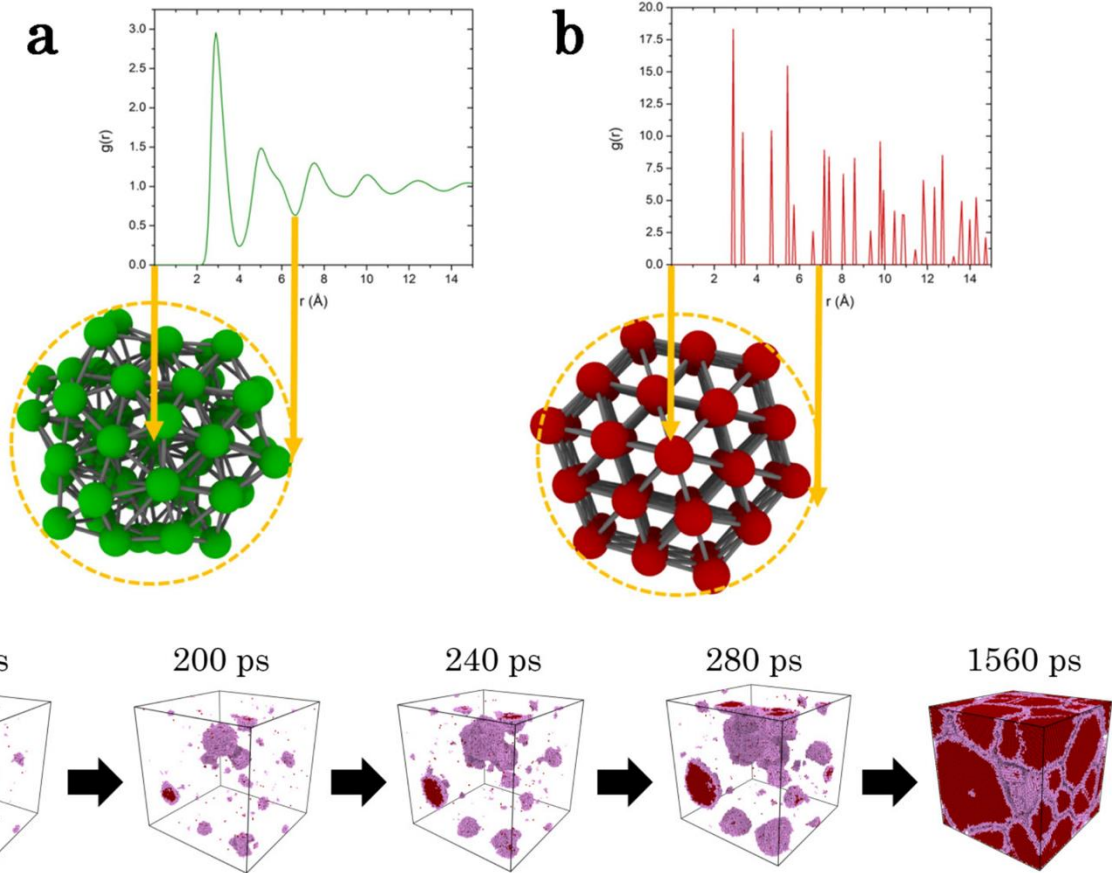*Universal fragment descriptors for predicting properties of inorganic crystals*

Xie *et al.*
**Phys Rev Lett (2018)**
*3-D Crystal Graph Convolutional Neural Networks for an Accurate and Interpretable Prediction of Material Properties*

# Interatomic potential models



$$\sum 7.51r^{3.98-3.93r}f(r)+\left(28.01+\sum -0.03r^{11.73-2.93r}f(r)\right)\left(\sum f(r)\right)^{-1}$$

**Mueller** *et al.*
**J Chem Phys (2020)**
*3-D Machine learning for interatomic potential models*

**Becker** *et al.*
**Sci report (2022)**
*Unsupervised topological learning approach of crystal nucleation*

33