# Unsupervised Learning for Chemistry:
# An introduction

Max Pinheiro Junior

max.pinheiro-jr@univ-amu.fr

SubNano

Aix*Marseille
université
Initiative d'excellence

"The study of algorithms that allow computer programs to **automatically improve through experience** or exposure to data."

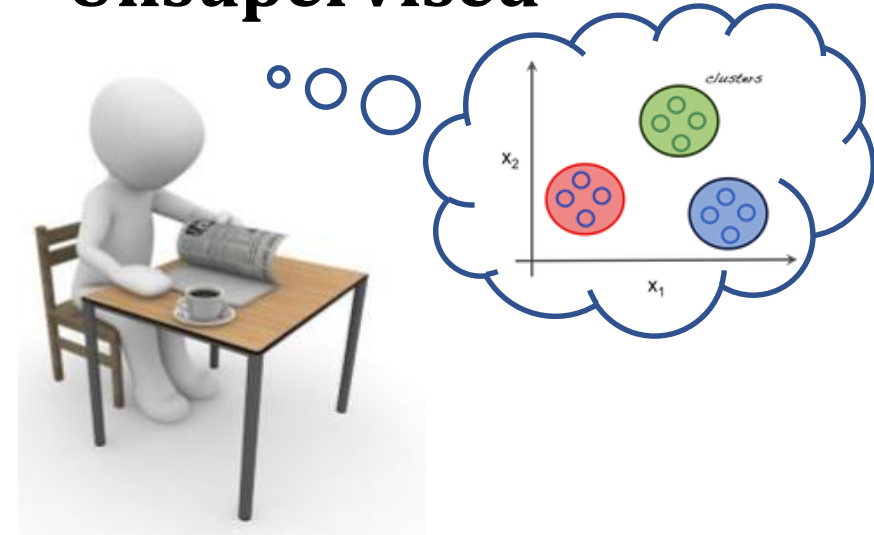(Tom Mitchel, Machine Learning, McGraw Hill, 1997)

## Supervised



- **Knowledge of output** (labelled data)
- **Goal:** Predict a specific quantity (class or value)
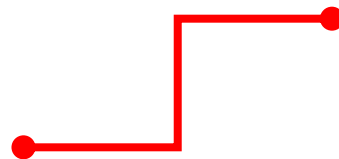- Can **measure accuracy** directly
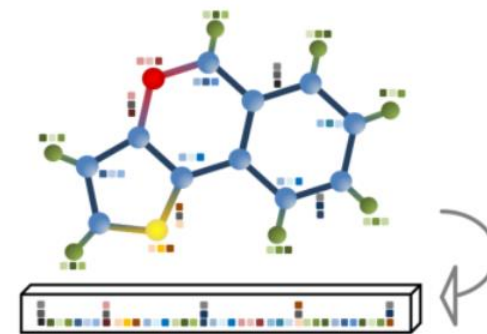
## Unsupervised



- **Unknown output** (unlabeled data)
- **Goal:** looking for structure or unusual patterns
- Indirect or qualitative evaluation
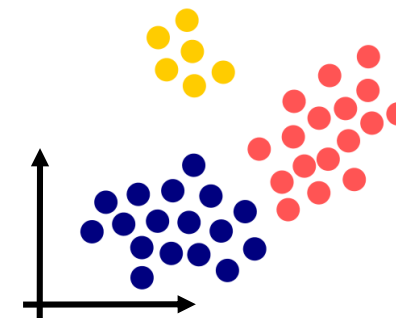
# Data-driven approach

1. Convert molecules into **numeric inputs**

2. Data preprocessing

3. Train or choose a **ML model**

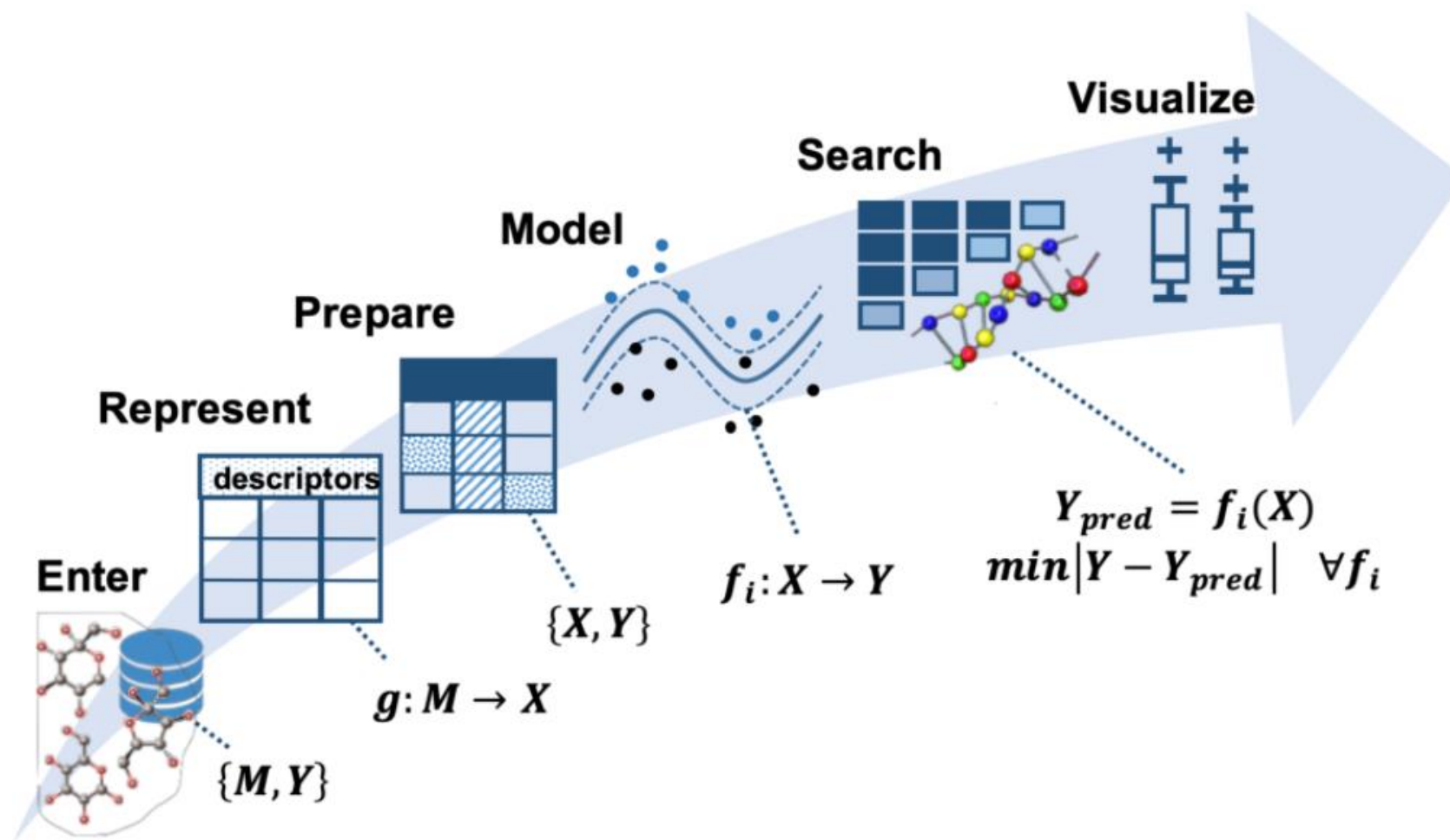4. **Evaluate** predictions or model outcomes
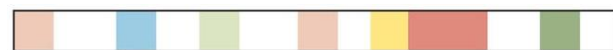
Representation

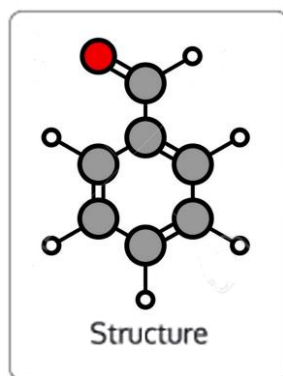Learning Model

Metrics for evaluation

**Fingerprint**

Structure

Descriptor

**Chemical properties**

# Molecular representations

## Fingerprint



Structure → Descriptor

↓

## Chemical properties

## Coulomb matrix

Simulates electrostatic interactions between atoms

$$M_{ij}^{\text{Coulomb}} = \begin{cases} 0.5 Z_i^{2.4} & \text{for } i = j \\ \dfrac{Z_i Z_j}{R_{ij}} & \text{for } i \neq j \end{cases}$$

$$
\begin{array}{c}
C \\ O \\ H \\ H \\ H \\ H
\end{array}
\begin{bmatrix}
36.9 & 33.7 & 5.5 & 3.1 & 5.5 & 5.5 \\
33.7 & 73.5 & 4.0 & 8.2 & 3.8 & 3.8 \\
5.5 & 4.0 & 0.5 & 0.35 & 0.56 & 0.56 \\
3.1 & 8.2 & 0.35 & 0.5 & 0.43 & 0.43 \\
5.5 & 3.8 & 0.56 & 0.43 & 0.5 & 0.56 \\
5.5 & 3.8 & 0.56 & 0.43 & 0.56 & 0.5
\end{bmatrix}
$$

M. Rupp *et al*, ***Phys. Rev. Lett.*** (2012) 108, 058301

General **Python** Libraries

**For Chemistry**

Data

Vis.

Math.

Modeling

Schnetpack, QML, …

# Part 1:
# Dimensionality Reduction

"

... dimensionality reduction yields a more compact, more easily interpretable **representation of the data**, focusing the attention on the **most relevant variables**.

"

— Page 289, Data Mining: Practical Machine Learning Tools and Techniques, 4th edition, 2016.

Main goal: decompose a dataset into a set of **orthogonal components** that explains a **maximum** amount of **variance**.



https://setosa.io/ev/principal-component-analysis/

PC as a new basis: $\mathbf{c}_{i1} = w_{11}x_{i1} + ... + w_{p1}x_{ip}$

Constrain (PCs normalized): $\sum_{j=1}^{p} w_{j1}^2 = 1$

**Objective**: Find w's that **maximize the variance** $\max_{w_{11},...,w_{p1}} \left\{ \frac{1}{n} \sum_{i=1}^{n} c_{i1}^2 \right\}$ **or**

**minimize the reconstruction error** $\min_{\mathbf{W}} \|\mathbf{X} - \mathbf{XCC^T}\|_F^2$ (information loss).



$$\mathbf{D}_3^2 = \mathbf{D}_1^2 + \mathbf{D}_2^2$$

| initial variance | = | remaining variance | + | lost variance |
|---|---|---|---|---|

$$\|\mathbf{a_i}\|^2 = \|w_i\,\mathbf{c}\|^2 + \|\mathbf{a_i} - w_i\,\mathbf{c}\|^2$$

| this is constant | | maximize this | **or** | minimize this |
|---|---|---|---|---|

PC as a new basis: $\mathbf{c}_{i1} = w_{11}x_{i1} + \ldots + w_{p1}x_{ip}$

Constrain (PCs normalized): $\displaystyle\sum_{j=1}^{p} w_{j1}^2 = 1$

**Objective**: Find w's that **maximize the variance** $\displaystyle\max_{w_{11},\ldots,w_{p1}} \left\{ \frac{1}{n}\sum_{i=1}^{n} \mathbf{c}_{i1}^2 \right\}$ **or**

**minimize the reconstruction error** $\displaystyle\min_{\mathbf{W}} \|\mathbf{X} - \mathbf{X}\mathbf{c}\mathbf{c^T}\|_F^2$ (information loss).

## Solution

Compute Eigenvalues and Eigenvectors of covariance matrix Σ:

$$\Sigma c_i = \lambda_i c_i \ (i = 1, 2, 3, \ldots, p), \qquad \Sigma = \frac{1}{n}\sum_{k=1}^{n} x_k^T x_k : \text{Covariance matrix}$$

Principal Component

# How to select the optimal number of components?

– Total variance = sum of variances of all individual principal components.

– The fraction of variance explained by a PC is the ratio between the variance of that PC and the total variance

$PC1 + PC2 + PC3 = 72\%$ **of total variance**



$$\frac{\sum_{i=1}^{3} \lambda_i}{\sum_{i=1}^{10} \lambda_i}$$

1. Compute the mean feature vector: $\mu = \dfrac{1}{p}\displaystyle\sum_{k=1}^{p} x_k$

2. Compute the covariance matrix: $\Sigma = \dfrac{1}{p}\displaystyle\sum_{k=1}^{p} (x_k - \mu)^T (x_k - \mu)$

3. Compute Eigenvalues and Eigenvectors of $\Sigma$:

   $\Sigma c_i = \lambda_i c_i \ (i = 1, 2, 3, \ldots, p), \qquad p =$ number of features

4. Select the $k$ eigenvectors $c_i$ corresponding to the largest eigenvalues $\lambda_i$

5. Reduce dimension from $p$ to $k$ with $Y = XC, \ C = (c_1, \ldots, c_k)$

Data with
isotropic variance

linear PCA

non-linear
transform.

**Kernel** Principal
Component Analysis

**Main idea:** map data onto a **high-dimensional feature space** (non-linear combinations) where data can be **linearly separable**

**Apply a non-linear transformation on the data:** $\mathbf{x} \to \Phi(\mathbf{x})$



$$\mathbf{x} = \begin{bmatrix} x_1 & x_2 \end{bmatrix}$$

$$\mathbf{x}' = \begin{bmatrix} x_1 & x_2 & x_1 x_2 & x_1^2 & x_2^2 & ... \end{bmatrix}$$

# Kernelized PCA

**Trick:** No need to know the mapping $\Phi(x)$ explicitly!

**Instead:** Use a **kernel (similarity) matrix** that can be expressed as dot products in the (HD) feature space.

$$K(x_i, x_j) = \langle \Phi(x_i), \Phi(x_j) \rangle = \begin{pmatrix} \Phi(x_1)\Phi(x_1) & \cdots & \Phi(x_1)\Phi(x_N) \\ \Phi(x_2)\Phi(x_1) & \cdots & \Phi(x_2)\Phi(x_N) \\ \vdots & \ddots & \vdots \\ \Phi(x_N)\Phi(x_1) & \cdots & \Phi(x_N)\Phi(x_N) \end{pmatrix}$$

**Examples of kernels:**

| Kernel function | Expression | Parameter |
|---|---|---|
| Polynomial kernel function | $K(x_i, x_j) = \left(x_i \cdot x_j + 1\right)^d$ | $d$ |
| Radial basis function (RBF) kernel function | $K(x_i, x_j) = \exp\left(-\gamma \lVert x_i - x_j \rVert^2\right)$ | $\gamma > 0$ |

1. Choose a kernel function – $K(x, x')$
    (ex: gaussian, polynomial, sigmoid, cosine)

2. Compute the kernel (similarity) matrix $K$ for all data points

3. Centralize the kernel matrix: $\widetilde{K} = K - \mathbf{1_N} K - K \mathbf{1_N} + \mathbf{1_N} K \mathbf{1_N}$

   <span style="color:red">Matrix with entries $1/N$</span>

4. Solve the eigenvalue equation: $\widetilde{K}\alpha_i = \lambda_i \alpha_i$

5. Select the $n$ largest eigenvalues to obtain the PCs in the feature space

$$y_j = \sum_{i=1}^{d} \alpha_{i,j} K(\mathbf{x}_i, \mathbf{x}) \text{ for } j = 1, 2, ..., N$$
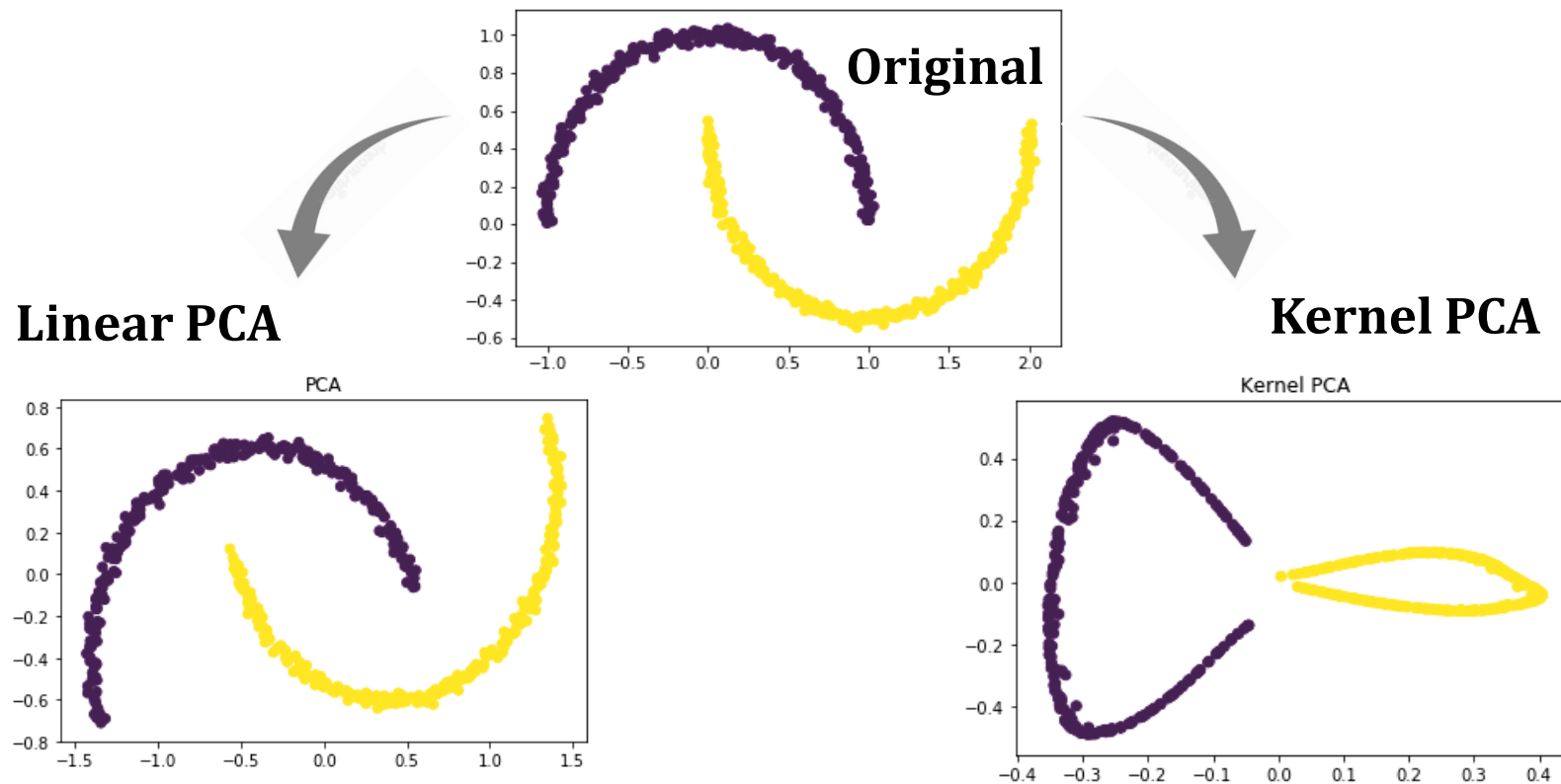
1. Choose a kernel function – $K(x, x')$
   (ex: gaussian, polynomial, sigmoid, cosine)

2. Compute the kernel (similarity) matrix $K$ for all data points

3. Centralize the kernel matrix: $\widetilde{K} = K - \mathbf{1_N}K - K\mathbf{1_N} + \mathbf{1_N}K\mathbf{1_N}$

4. Solve the eigenvalue equation: $\widetilde{K}\alpha_i = \lambda_i\alpha_i$

5. Select the $n$ largest eigenvalues to obtain the PCs in the feature space

$$y_j = \sum_{i=1}^{d} \alpha_{i,j}K(\mathbf{x}_i, \mathbf{x}) \text{ for } j = 1, 2, ..., N$$

**Original**

**Linear PCA**

**Kernel PCA**

```python
from sklearn.decomposition import PCA
pca = PCA(n_components = 2)
X_pca = pca.fit_transform(X)

plt.title("PCA")
plt.scatter(X_pca[:, 0], X_pca[:, 1], c = y)
plt.show()
```
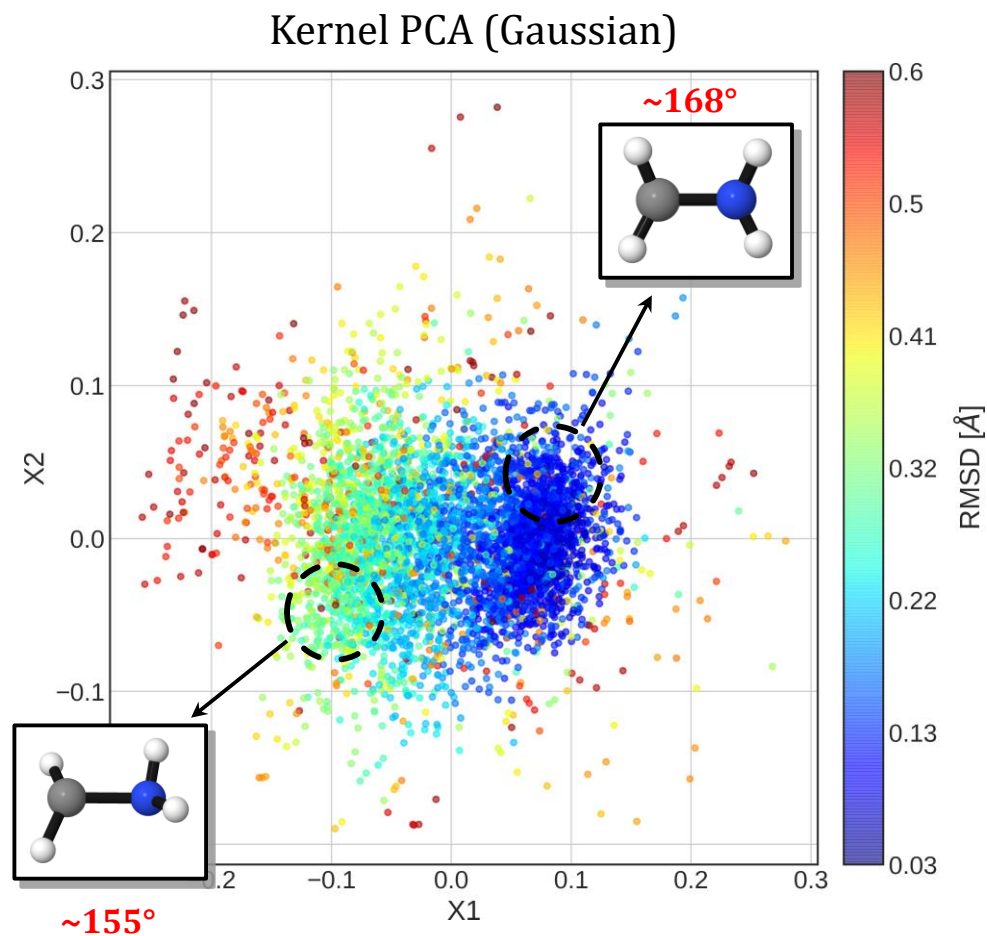
```python
from sklearn.decomposition import KernelPCA
kpca = KernelPCA(kernel ='rbf', gamma = 15)
X_kpca = kpca.fit_transform(X)

plt.title("Kernel PCA")
plt.scatter(X_kpca[:, 0], X_kpca[:, 1], c = y)
plt.show()
```

Kernel PCA (Gaussian)

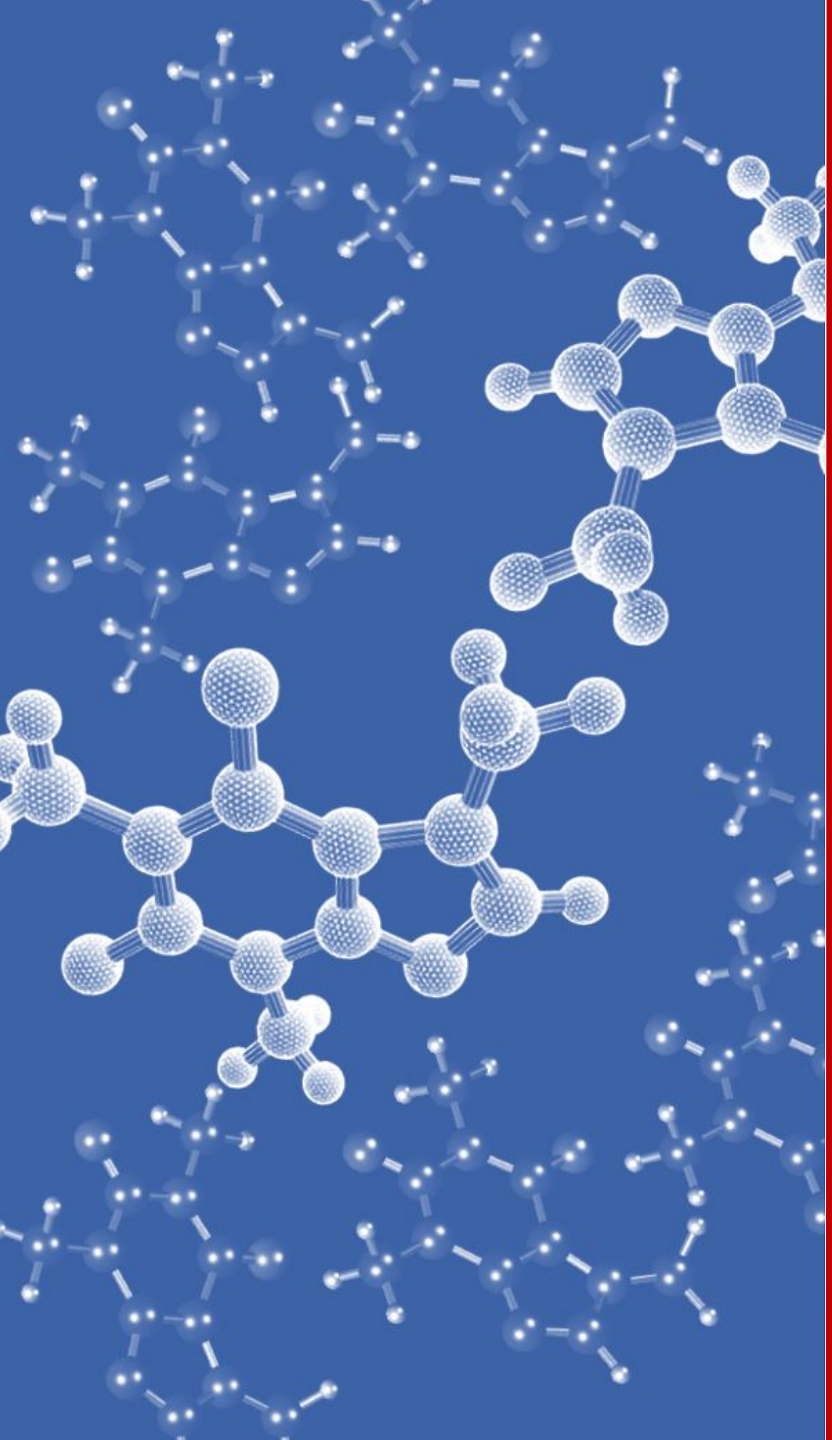**Most of geometry changes occur along $X1$**

**Blue cluster:** similar to the S0 equilibrium geometry (small RMSD)

**Cyan cluster:** rotated NH2 group (near conical intersection?)

**Red points:** highly distorted or broken geometries

**K-PCA:**

- Take nonlinearities into account
- Identify reaction paths

# Part 2:
# Clustering Methods

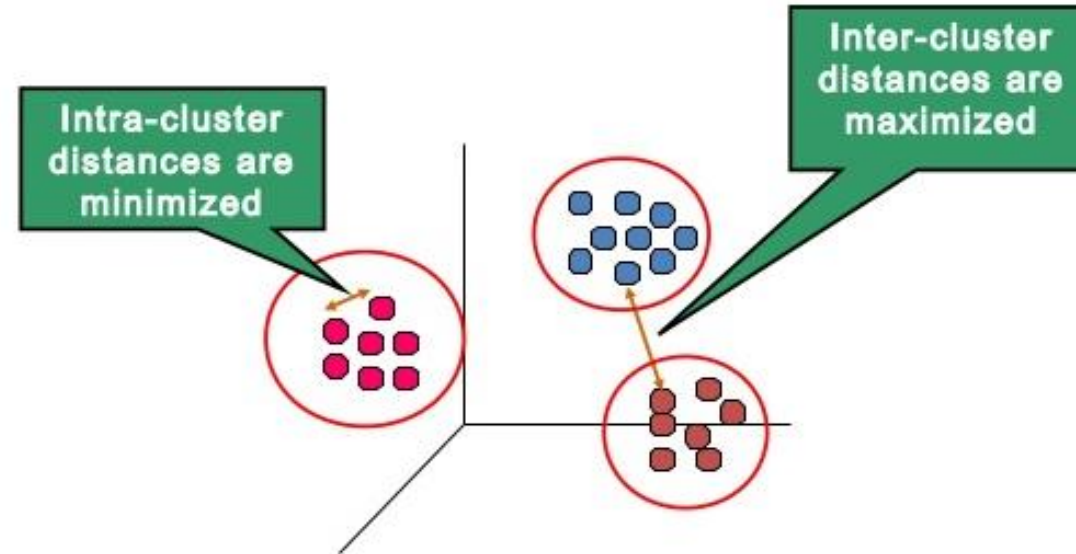How to **find patterns in data** with unknown labels?



**Task: group data by similarity**

> Finding groups of objects such that the objects within a group **be similar (or related)** to one another and **different from (or unrelated to)** the objects in other groups.
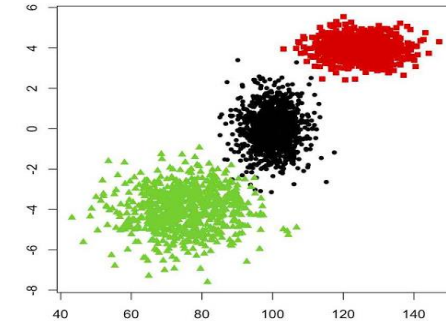
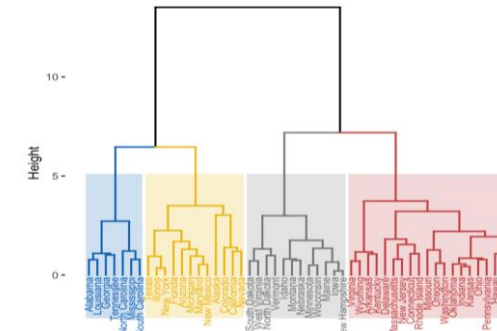— Page 490, Tan et al, Introduction to Data Mining, 1st edition, 2006.

**1. Partitional (K-Means**, K-Medoids**):**
divides data into **non-overlapping** groups based on distances. Useful for **spherical shape** clusters.
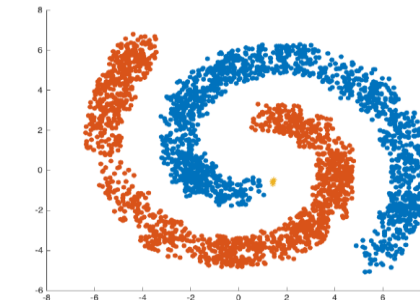
**2. Hierarchical (Single-linkage):**
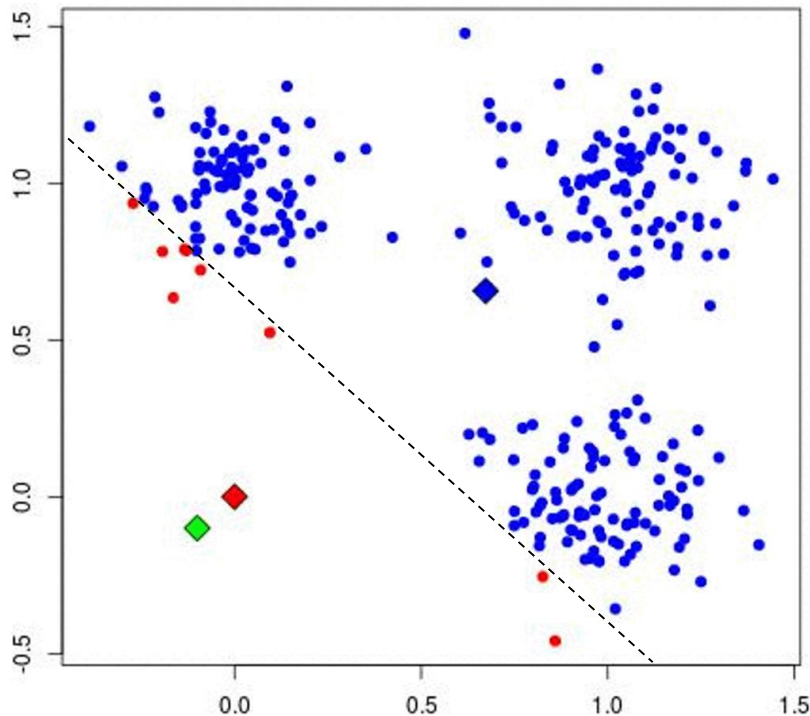determines cluster assignments by building a hierarchy.

**3. Density-based (DBSCAN):**
Clusters are assigned to dense regions of data space. Can find arbitrarily shaped clusters. May filter out outliers.

# K-Means clustering

**Main idea:** partition the data into **K clusters** by assigning each observation to the group of points with the **nearest mean** (cluster centroid).



**Algorithm**

**Input:** data points $\mathbf{X} = [x_1, x_2, \ldots, x_n], x_i \in \mathbb{R}^n$

**Parameter:** number of clusters $K$

**Output:** centroids of the $K$ clusters

1. Initialize the $K$ cluster centers;

**while** *centroids change* **do**

   2. Calculate distances between the centroids and **all data points**;

   3. Assign each data point to its **nearest centroid**;

   4. **Update the centroids** using the mean of the of the current cluster memberships;

**end**

https://www.naftaliharris.com/blog/visualizing-k-means-clustering/

# K-Means clustering

**Main idea:** partition the data into **K clusters** by assigning each observation to the group of points with the **nearest mean** (cluster centroid).



**Algorithm**

**Input:** data points $\mathbf{X} = [x_1, x_2, ..., x_n], x_i \in \mathbb{R}^n$

**Parameter:** number of clusters $K$

**Output:** centroids of the $K$ clusters

1. Initialize the $K$ cluster centers;

while *centroids change* **do**

   2. Calculate distances between the centroids and **all data points**;

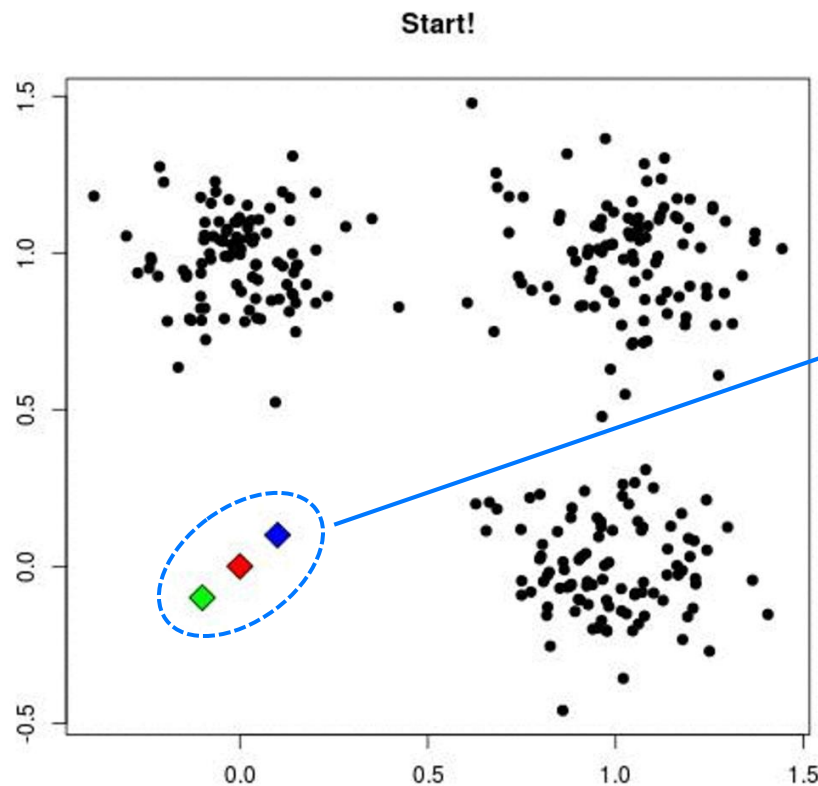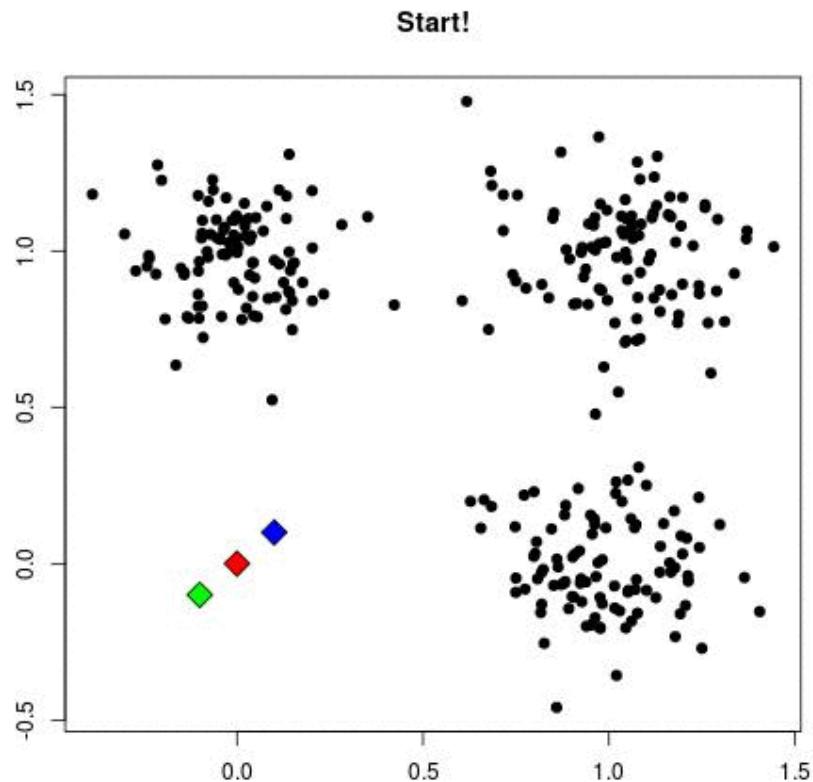   3. Assign each data point to its **nearest centroid**;

   4. **Update the centroids** using the mean of the of the current cluster memberships;

**end**

# K-Means clustering

**Main idea:** partition the data into **K clusters** by assigning each observation to the group of points with the **nearest mean** (cluster centroid).



Start!

**Algorithm**

**Input:** data points $\mathbf{X} = [x_1, x_2, ..., x_n], x_i \in \mathbb{R}^n$

**Parameter:** number of clusters $K$

**Output:** centroids of the $K$ clusters

1. Initialize the $K$ cluster centers;

**while** *centroids change* **do**

   2. Calculate distances between the centroids and **all data points**;

   3. Assign each data point to its **nearest centroid**;

   4. **Update the centroids** using the mean of the of the current cluster memberships;

**end**

https://www.naftaliharris.com/blog/visualizing-k-means-clustering/

The **objective function** to be minimized is the Sum of Squared Error (SSE), given by:

$$J = \sum_n \sum_k r_{nk} \|x_n - \mu_k\|^2$$

**to optimize**

Where $\begin{cases} x \text{ is a sample } n, \ x_n \in \mathbb{R}^D \\ \mu_k \text{ is the cluster center of cluster } k, \mu_k \in \mathbb{R}^D \\ \boldsymbol{r} \text{ is the cluster membership}, \mathbf{r} \in \mathbb{R}^{N \times D} \end{cases}$

$$r_{nk} = \begin{cases} 1, & \text{if } x_n \in \text{cluster } k \\ 0, & otherwise \end{cases} \longrightarrow \arg\min_{k'} \left\| x_n - \mu_{k'} \right\|^2 = k$$

$$\mu_k^* = \arg\min_{\mu_k} J \longrightarrow \text{optimal cluster center}$$

$$\sum_n r_{nk}\mu_k - \sum_n r_{nk}x_n = 0$$

$$\frac{\delta J}{\delta \mu_k} = \frac{\delta\left[\sum_n \sum_k r_{nk}\|x_n - \mu_k\|^2\right]}{\delta_{\mu_k}} = 0$$

$$\mu_k \sum_n r_{nk} - \sum_n r_{nk}x_n = 0$$

$$\sum_n r_{nk} \times 2(x_n - \mu_k)(-1) = 0$$

$$\boxed{\mu_k = \frac{1}{\sum_n r_{nk}} \sum_n r_{nk}x_n}$$

$$\sum_n 2r_{nk}\mu_k - \sum_n 2r_{nk}x_n = 0$$

**Centroid definition**

# K-Means: *Pros* and *Cons*

Simple, intuitive and easy to implement

K-means works well when the data form **compact clouds** with globular shapes, and **well separated** from one another

Number of clusters should be provided by the user ($K = ?$)

└─────▶ **Solution:** elbow method

**Sensitive to outlier** points, which can affect the mean values significantly

└─────▶ **Solution:** data preprocessing!

**Sensitive to initialization:** clustering results may vary significantly with the initial guess for the cluster centers

└─────▶ **Solution:** multiple executions with random initializations

# K-Means in molecular simulations



**2020**

Overcoming the Heuristic Nature of *k*-Means Clustering: Identification and Characterization of Binding Modes from Simulations of Molecular Recognition Complexes

Parker Ladd Bremer, Danna De Boer, Walter Alvarado, Xavier Martinez, and Eric J. Sorin*

How to **automatically identify** different types of **protein-inhibitor complexes** from **MD data**? Ans: **K-Means!**



**K = 4**

**K = 3**

✅ Unsupervised learning is quite challenging
   - input = data + some prior knowledge + intuition.

✅ Unsupervised ≈ exploratory data analysis → **try to feel the different "flavors" of your data!**

✅ Unsupervised can be useful for labeled data. Validate your analysis!

✅ Dimensionality reduction: projected data does not always reflect the relationships/distance of the original data space.

✅ Data quality for clustering:
   - choose an adequate representation for the chemical system
   - invest time in preprocessing
   - which distance metric should I use?

**Tutorial in Python:** github.com/maxjr82/CECAM-MLQCDyn

## Analysing dynamics data with unsupervised learning

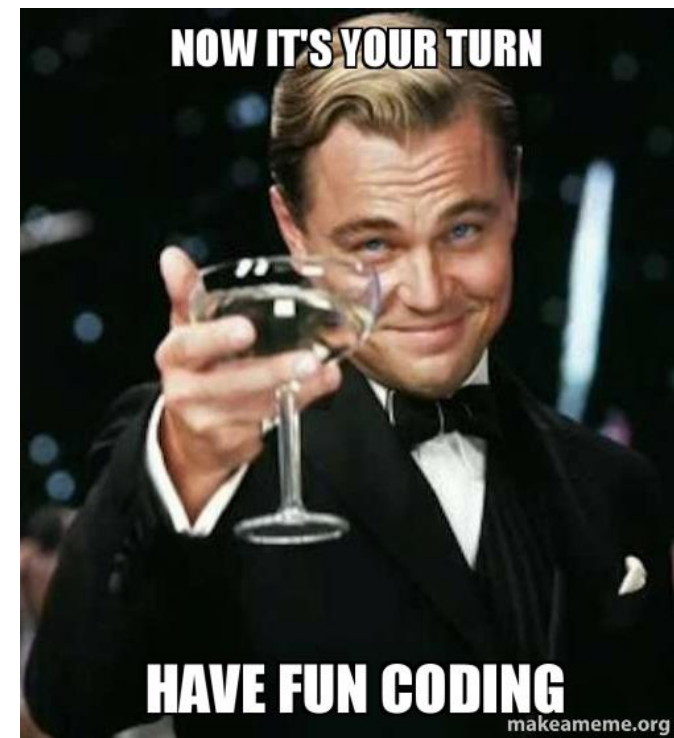### Machine Learning and Quantum Computing for Quantum Molecular Dynamics - 2022

Author: Max Pinheiro Jr[1]

[1]Aix Marseille University, CNRS, Marseille, France
[Last updated: August 24, 2022]

Aix*Marseille
université
*initiative d'excellence*

**::cecam**
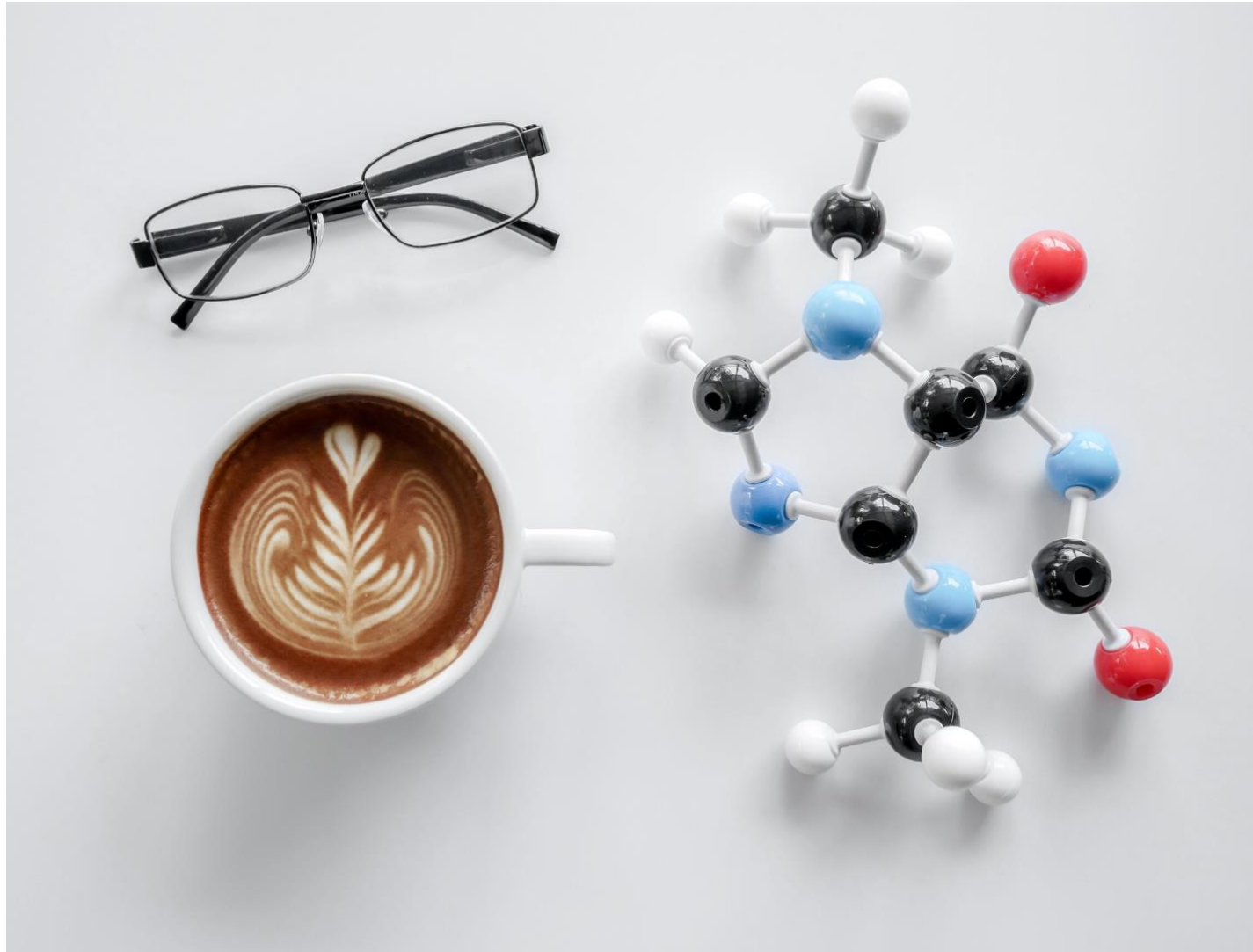Centre Européen de Calcul Atomique et Moléculaire

This tutorial is devoted to the application of unsupervised machine learning methods to analyze and interpret chemical data generated by molecular dynamics simulations. Broadly speaking, the main goal in unsupervised learning is to find natural grouping structure or possible associations within the data, or even a compact representation for the data based on measures on a set of inputs. Many algorithms have been developed to accomplish these tasks being *dimensionality reduction* and *clustering analysis* the two most representative and important branches of unsupervised learning. Here we will explore some of the these techniques for data analysis by considering practical examples based on a molecular configurational dataset generated from snapshots of molecular dynamics simulations to facilitate your understanding of the theory underpinning the methods. Instead of going deep into mathematical details we will dive into a few introductory examples of the methods (when viable) using simplified or *toy* data to explain and illustrate how the algorithms work in practice.

The tutorial was designed to run in a Python environment, so a basic knowledge of this programming language will be helpful. For those that are not familiar with Python, I suggest you to take a time to check the Python documentation and the Numpy manual for clues.

**GitHub**      **Google Colab**

NOW IT'S YOUR TURN
HAVE FUN CODING
makeameme.org

**Send questions to:**
**maxjr82@gmail.com**

# Thank you!