

# Lecture 12

## Fuzzy clustering by PCCA+

### CONTENTS

A. Crisp clustering	1
B. Fuzzy clustering	2
C. Robust Perron Cluster Cluster Analysis	2
D. Solution for $n_c = 2$	4
E. Example: triple-well potential	7
F. Example: periodic triple-well potential	8
References	9

### A. Crisp clustering

Clustering or coarse-graining is a technique used in data analysis to group similar data points or objects together based on certain features or characteristics they share. In the context of dynamical systems (deterministic and stochastic), clustering refers to the discretization of the state space into subsets containing states with similar static and kinetic properties (e.g. equilibrium distribution and rates). Clustering is useful because it allows to represent continuous dynamic as a discrete process. Consequently, continuous objects such as operators and functions can be represented by discrete objects such as matrices and vectors that are easier to use in practical applications.

For example, consider a dynamical system defined on the state space  $\Gamma$  discretized with a Voronoi tessellation of  $k$  disjoint cells, or clusters,  $\Gamma_i$  such that  $\Gamma = \cup_i^k \Gamma_i$ , where each cell  $\Gamma_i$  is defined by the indicator function

$$\mathbf{1}_i(x) = \begin{cases} 1 & \text{if } x \in \Gamma_i , \\ 0 & \text{if } x \notin \Gamma_i . \end{cases} \quad (1)$$

The choice of the tessellation is arbitrary, it could be a tessellation made of either regular or irregular  $N_d$ -polytopes (polygons in 2D, polyhedra in 3D).

This kind of clustering is known as crisp clustering and is largely used. However, it has numerous limitations:

- Although dynamics is Markovian, its cluster representation may lose this property.
- Crisp clustering is not robust to noise. Small perturbations in the dynamics could be amplified in its cluster representation.
- Crisp clustering methods struggle with identifying clusters of irregular shapes or clusters that are connected but separated by sparse regions.
- When analysing metastable regions, it is not always possible to uniquely determine the boundaries of metastable regions.
- Crisp clustering requires the specification of the number of clusters  $k$  a priori. Increasing the resolution, i.e. the number of clusters, may alleviate some problems, but the initial dataset may not have enough data points. Furthermore, a high number of clusters may require more resources to perform calculations with the matrices involved.

To address some of these limitations, researchers often explore alternative clustering approaches, such as fuzzy clustering, hierarchical clustering, and density-based clustering, which offer more flexibility and improved performance in specific scenarios.

## B. Fuzzy clustering

Fuzzy clustering, in contrast to crisp clustering, allows for the assignment of data points to multiple clusters with varying degrees of membership. Instead of assigning each data point to a single cluster (as in crisp clustering), fuzzy clustering assigns a membership value to each data point for each cluster, indicating the degree (or the probability) to which the point belongs to that cluster. This provides a more nuanced representation of the inherent uncertainty or ambiguity in the data.

## C. Robust Perron Cluster Cluster Analysis

We use fuzzy clustering to identify the  $n_c$  metastable states also referred to as metastable macro-state, or conformations, of a system driven by Langevin dynamics, with a potential

energy function  $V(x) : \Gamma \subset \mathbb{R}^{N_d} \rightarrow \mathbb{R}$  with  $n_c$  minima separated by energy barriers higher than the thermal energy  $k_B T$ . For example, given a double well potential:  $n_c = 2$ ; for a triple-well potential  $n_c = 3$ ; ....

For each  $i$ th macro-state, we introduce the membership function

$$\chi_i(x) \in [0, 1], \quad (2)$$

also known as almost characteristic function, that indicates the probability, or membership degree, that a state  $x$  belongs to the cluster  $j$ . The membership functions fulfil the partition of the unit

$$\sum_i^{n_c} \chi_i(x) = 1. \quad (3)$$

Note that, if we indicate  $\chi$  without the sub-index  $i$ , we refer to the set  $\chi = \{\chi_1, \chi_2, \dots, \chi_{n_c}\}$  containing all the  $n_c$  membership functions organized in columns. From a geometrical point of view, the points of the  $\chi$  functions lie on the standard  $(n_c - 1)$ -simplex as illustrated in fig. 1. The standard term indicates that the vertices of the simplex are the unit vectors  $e_1, e_2, \dots, e_{n_c}$ .

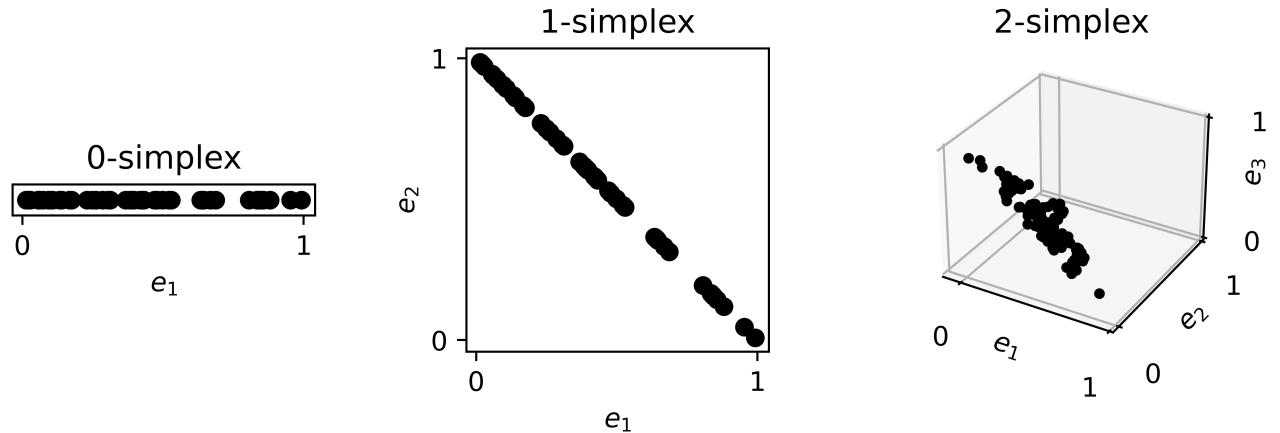


FIG. 1. Random points that fulfill the partition of unity.

The robust Perron Cluster Cluster Analysis (PCCA+) algorithm [1–4] determines the membership functions as a linear combination of the first  $n_c$  dominant eigenfunctions  $\psi = \{\psi_0, \psi_1, \dots, \psi_{n_c}\}$  of the infinitesimal generator

$$\mathbf{Q}\psi_i = \kappa_i\psi_i, \quad (4)$$

or the associated Koopman operator

$$\mathcal{K}_\tau \psi_i = \lambda_i(\tau) \psi_i. \quad (5)$$

We recall that the first eigenfunction  $\psi_0$  is equal to 1, while the other eigenfunctions have positive and negative values, and represent the dominant processes that constitute the dynamics of the system.

Similar to the  $\chi$  functions, also the dominant  $n_c$  eigenfunctions form an  $(n_c - 1)$ -simplex, the vertices of which, however, are not the unit vectors (see figs. E and F as examples). The idea underlying PCCA+ is then to find the linear transformation such that

$$\chi = \psi A, \quad (6)$$

where  $A$  is a matrix of size  $n_c \times n_c$ . The simplex has a physical interpretation: the vertices represent the conformations of the system, the points on the edges represent the transition regions. Additionally, the membership functions allow the direct Galerkin discretization of the infinitesimal generator

$$\mathbf{Q}_c = \langle \chi, \chi \rangle_\pi^{-1} \langle \chi, \mathcal{Q}\chi \rangle_\pi, \quad (7)$$

where  $\mathbf{Q}_c$  is an  $n_c \times n_c$  matrix whose entries express the transition rates between fuzzy sets.

#### D. Solution for $n_c = 2$

Unfortunately, determining the matrix  $A$  is not easy, as there are an infinite number of possible solutions, which can only be determined solving an optimization problem after appropriate objective functions have been defined [1]. However, for the sole case when  $n_c = 2$ , a unique solution can be determined.

If  $n_c = 2$  the matrix  $A$  reads.

$$A = \begin{cases} a_{00} & a_{01} \\ a_{10} & a_{11} \end{cases}. \quad (8)$$

First, we pose the following constraints on  $A$ :

1.

$$\chi = \psi A \rightarrow \chi_i(x) = \sum_{j=0}^{n_c} a_{ji} \psi_j(x),$$

2.

$$\sum_{i=0}^{n_c-1} \chi_i(x) = 1,$$

3.

$$\chi_i(x) \leq 0.$$

Then, we rewrite the first condition as

$$\chi_i(x) = \sum_{j=0}^{n_c-1} a_{ji} \psi_j(x) \quad (9)$$

$$= a_{0i} \psi_0 + \sum_{j=1}^{n_c-1} a_{ji} \psi_j(x), \quad (10)$$

and applying the second condition, we obtain an expression for  $a_{0i}$ :

$$a_{0i} \psi_0(x) + \sum_{j=1}^{n_c-1} a_{ji} \psi_j(x) \leq 0 \quad (11)$$

$$a_{0i} \psi_0(x) \leq - \sum_{j=1}^{n_c-1} a_{ji} \psi_j(x) \quad (12)$$

$$a_{0i} = - \min_x \sum_{j=1}^{n_c-1} a_{ji} \psi_j(x), \quad (13)$$

where we used  $\psi_0(x) = 1$ . We now use the third condition to rewrite the second as

$$\sum_{i=0}^{n_c-1} \chi_i(x) = \quad (14)$$

$$\sum_{i=0}^{n_c-1} \sum_{j=0}^{n_c} a_{ji} \psi_j(x) = \quad (15)$$

$$\sum_{j=0}^{n_c-1} \sum_{i=0}^{n_c} a_{ji} \psi_j(x) = \quad (16)$$

$$\sum_{j=0}^{n_c-1} \delta_{j0} \psi_j(x) = 1. \quad (17)$$

where  $\delta_{j0}$  is the Kronecker-delta and we used again  $\psi_0(x) = 1$ . From the equality

$$\sum_{i=0}^{n_c} a_{ji} \psi_j(x) = \delta_{j0} \quad (18)$$

$$a_{j0} + \sum_{i=1}^{n_c} a_{ji} \psi_j(x) = \delta_{j0}, \quad (19)$$

then we obtain an expression for  $a_{j0}$ :

$$a_{j0} = \delta_{j0} - \sum_{i=1}^{n_c} a_{ji} \psi_j(x), \quad (20)$$

Applying  $i = 0, 1$  and  $j = 0, 1$  to eqs. 13, 20 yields:

$$\begin{cases} a_{00} = -\min_x a_{10} \psi_1(x), \\ a_{01} = -\min_x a_{11} \psi_1(x), \\ a_{00} = 1 - a_{01}, \\ a_{11} = -a_{10}. \end{cases} \quad (21)$$

From the forth equality, the second one becomes

$$a_{00} = a_{11} \max_x \psi_1(x). \quad (22)$$

Finally we have two equations for  $\psi_1$ :

$$\max_x \psi_1(x) = \frac{a_{00}}{a_{11}}, \quad (23)$$

and

$$-\min_x \psi_1(x) = \frac{a_{01}}{a_{11}}. \quad (24)$$

Their sum yields

$$\max_x \psi_1(x) - \min_x \psi_1(x) = \frac{a_{00} - a_{01}}{a_{11}} = \frac{1}{a_{11}}. \quad (25)$$

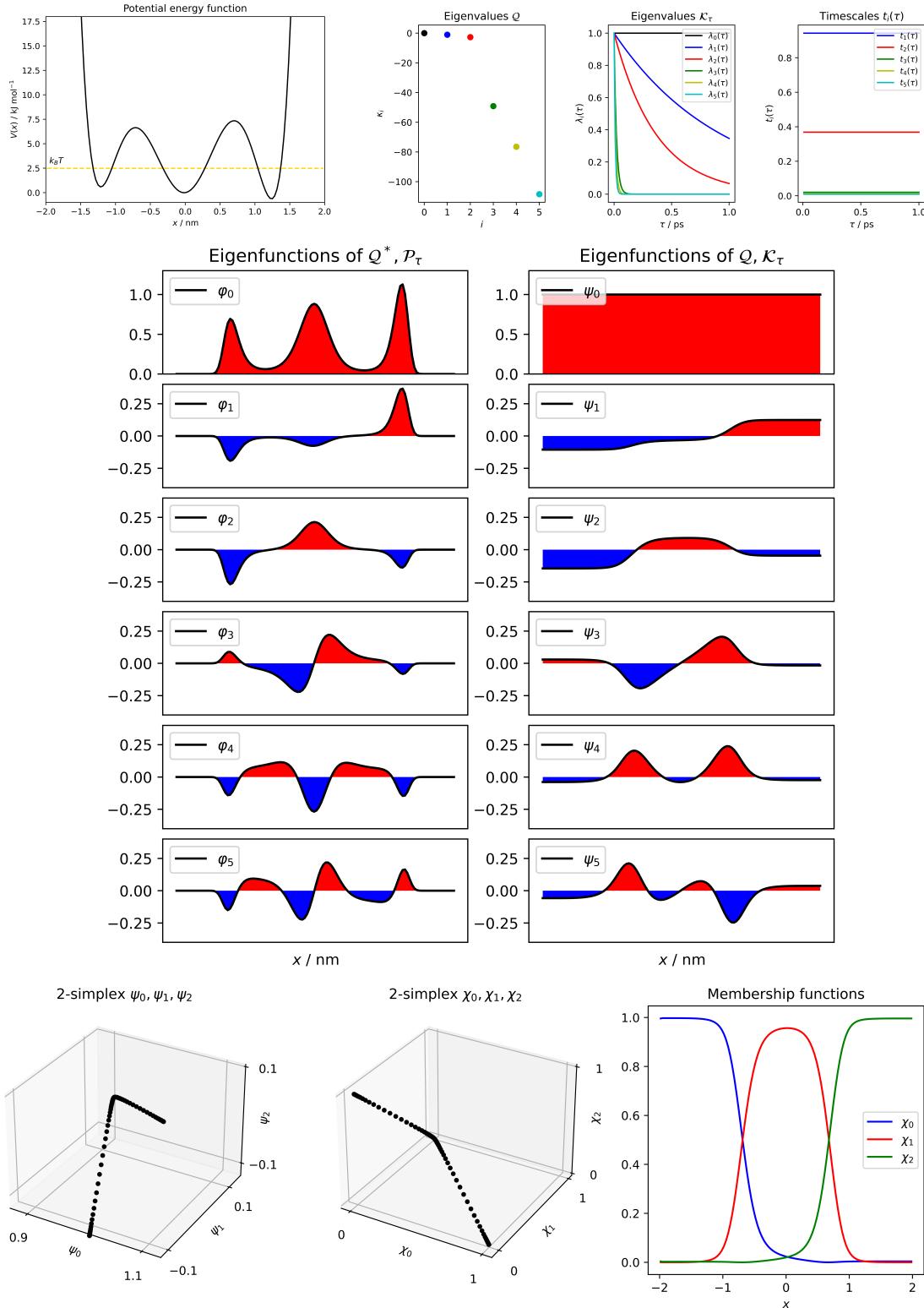
Finally, one obtains an expression for each entry of the matrix  $A$ :

$$\begin{cases} a_{00} = \frac{\max_x \psi_1(x)}{\max_x \psi_1(x) - \min_x \psi_1(x)} \\ a_{01} = -\frac{\min_x \psi_1(x)}{\max_x \psi_1(x) - \min_x \psi_1(x)} \\ a_{10} = -\frac{1}{\max_x \psi_1(x) - \min_x \psi_1(x)} \\ a_{11} = \frac{1}{\max_x \psi_1(x) - \min_x \psi_1(x)}, \end{cases} \quad (26)$$

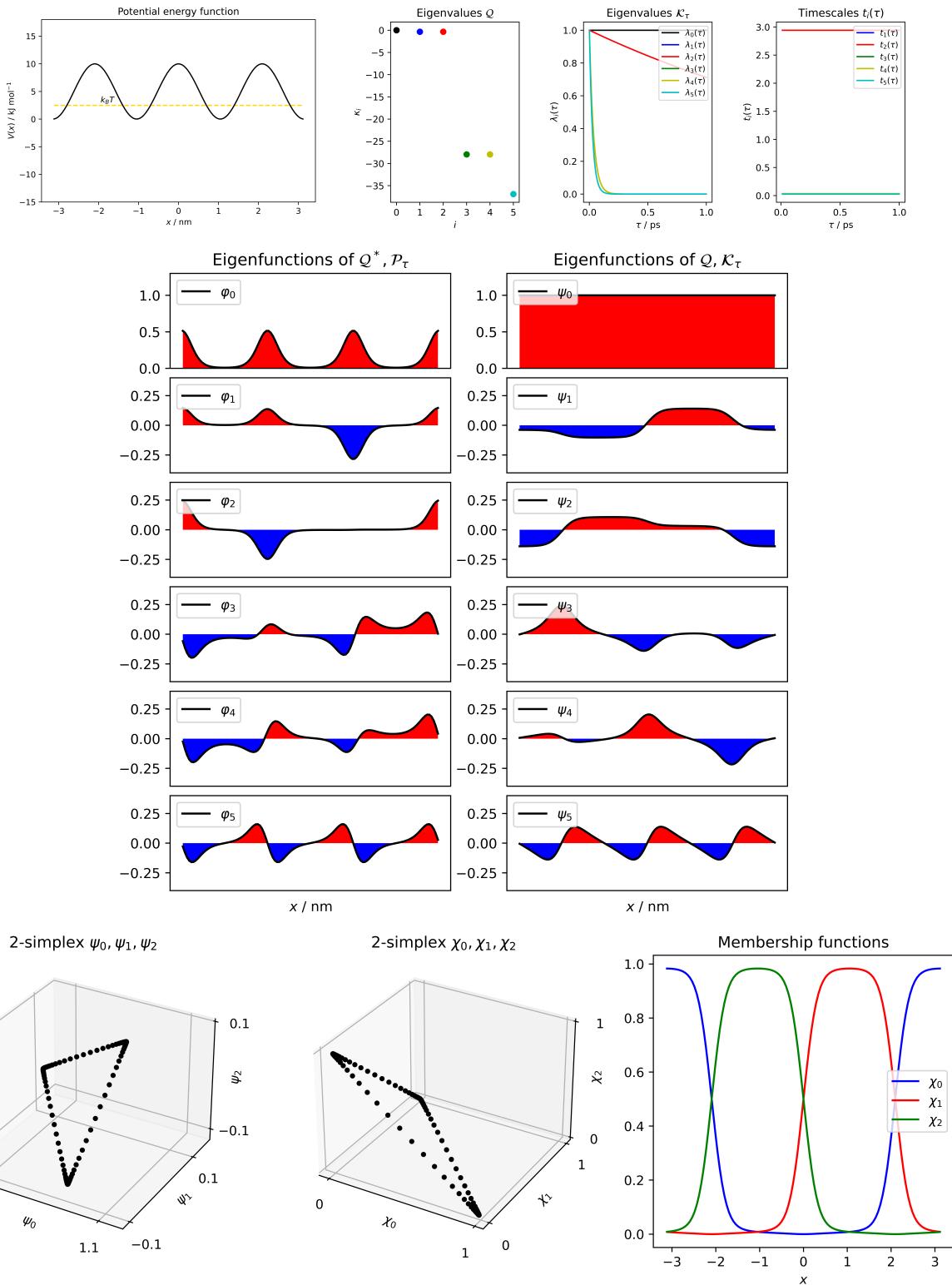
The two membership functions for the case  $n_c = 2$  read

$$\begin{cases} \chi_1(x) = \frac{\max_x \psi_1(x) - \psi_1}{\max_x \psi_1(x) - \min_x \psi_1(x)}, \\ \chi_2(x) = \frac{\psi_1 - \min_x \psi_1(x)}{\max_x \psi_1(x) - \min_x \psi_1(x)} = 1 - \chi_1(x). \end{cases} \quad (27)$$

### E. Example: triple-well potential



### F. Example: periodic triple-well potential



- [1] P. Deuflhard and M. Weber, Robust perron cluster analysis in conformation dynamics, *Linear Algebra Appl.* **398**, 161 (2004).
- [2] S. Kube and M. Weber, A coarse graining method for the identification of transition rates between molecular conformations, *J. Chem. Phys.* **126**, 024103 (2007).
- [3] S. Röblitz and M. Weber, Fuzzy spectral clustering by pcca+: Application to markov state models and data classification, *Adv. Data Anal. Classif.* **7** (2013).
- [4] M. Weber, Implications of pcca+ in molecular simulation, *Computation* **6** (2018).