

Clustering Methods

So far ... **Supervised Learning**

$x_1, x_2 \dots x_n$ inputs

$y_1, y_2 \dots y_n$ outputs.

New input x : Predict y .

Today : **Unsupervised Learning.**

$x_1, x_2 \dots x_n$ inputs

No Outputs.

Goal : identify relevant subgroups,
called **clusters**.

Example: images car / planes. No labels

identify: Cluster 1: car

Cluster 2: planes

① k-means

$$X = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$$

Goal of k-means:

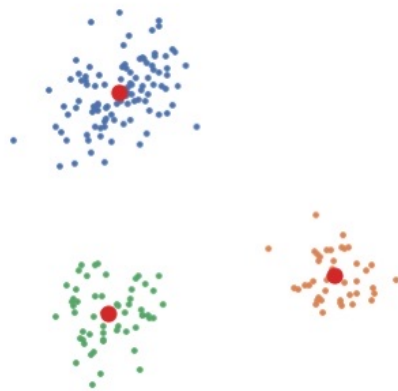
Find a set of k representatives,
called centroids.

$k=3$

● centroids

● observations

Create the clusters
by grouping points
the closest to the same
centroid.



How should the centroids be defined?

• $k=1$: $X = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$

The best way to summarize X using one point:

→ Choose $y^* = \frac{x_1 + \dots + x_n}{n}$ the mean.

Actually, y^* is the minimizer of

$$F: y \in \mathbb{R}^d \mapsto \frac{1}{n} \sum_{i=1}^n \|x_i - y\|^2.$$

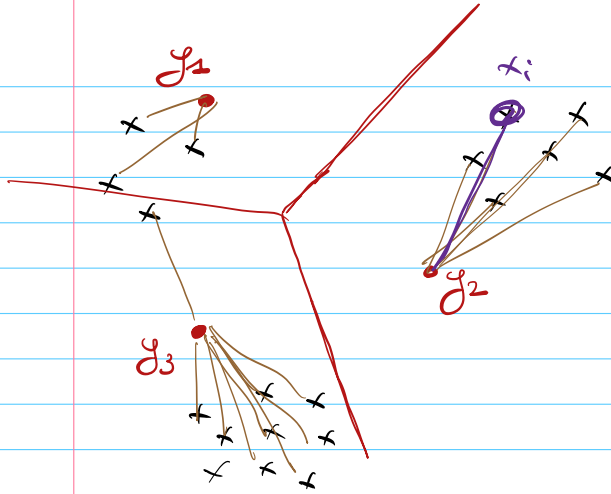
proof: $\nabla F(y) = \frac{1}{n} \sum_{i=1}^n 2(y - x_i)$

y such that $\nabla F(y) = 0 = \frac{1}{n} \sum_{i=1}^n (y - x_i)$

→ $y = y^* = \frac{1}{n} \sum_{i=1}^n x_i$.

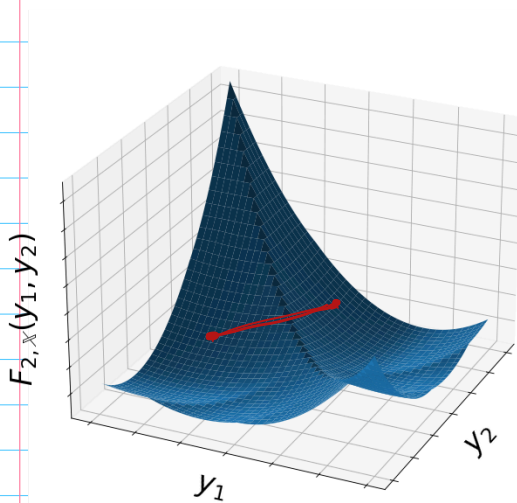
For $k > 1$: Generalization

$$F_{k,X}: (y_1, \dots, y_k) \in (\mathbb{R}^d)^k \mapsto \frac{1}{n} \sum_{i=1}^n \min_{l=1, \dots, k} \|y_l - x_i\|^2$$



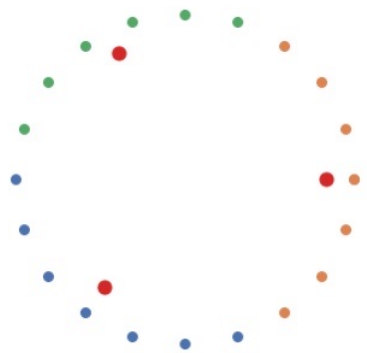
Choose (y_1, y_2, y_3)
to minimize
 $\sum (-)^2$

\Rightarrow The k -means of X is the ¹¹⁹ minimizer (y_1^*, \dots, y_k^*) of $F_{k,X}$.
 \hookrightarrow also called centroids



The function $F_{k,X}$
is NOT convex.
($k \geq 2$)

$k=2$ $d=1$



Several k -means
may exist.

How to compute the k -means?

LLOYD'S ALGORITHM

Initialization: centroids y_1^0, \dots, y_k^0

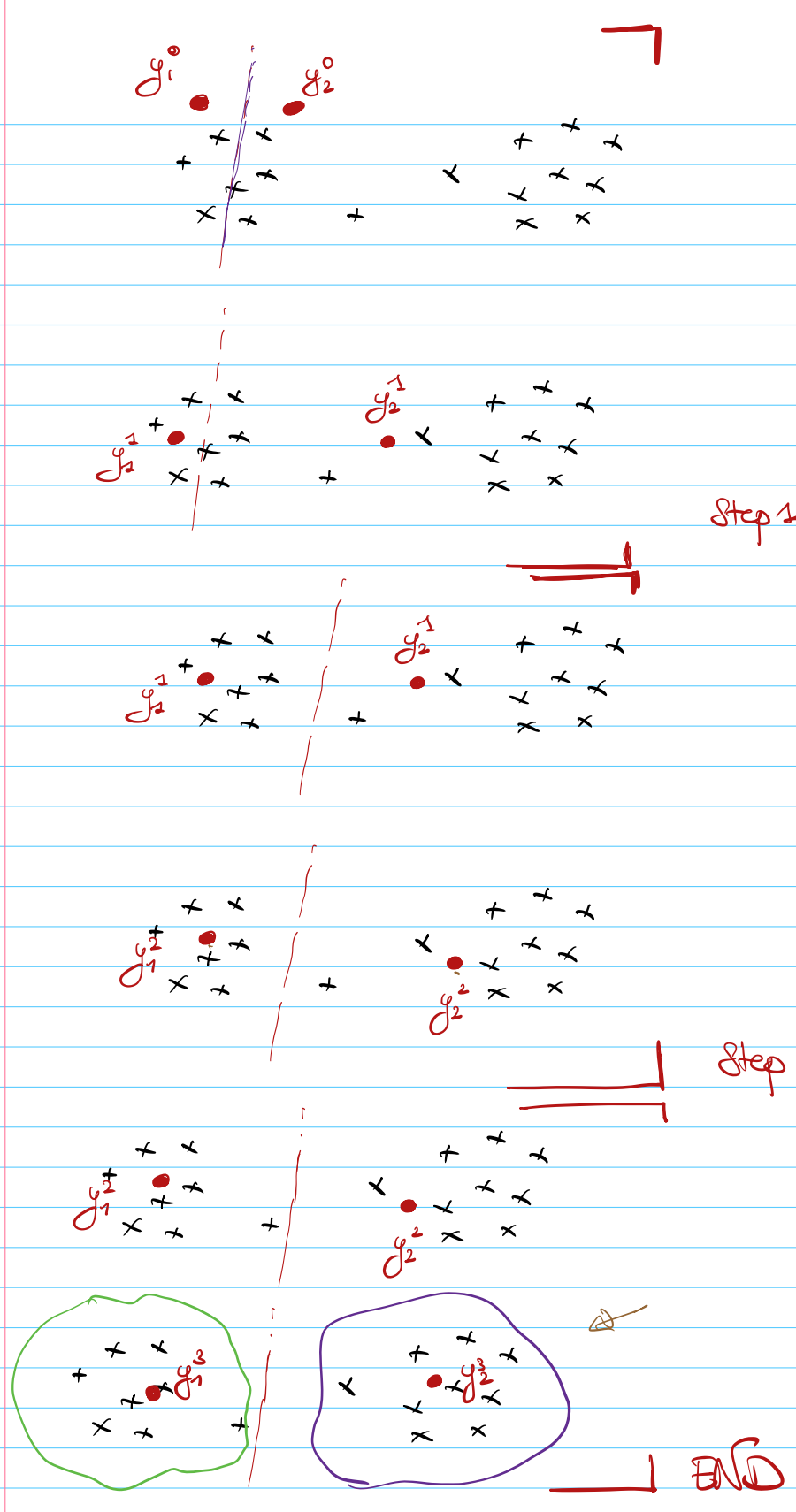
For $t=0 \dots T-1$:

$\forall l=1 \dots k$:

cluster l at
time t

$$\left[\begin{array}{l} I_l^t = \{ i : y_l^t \text{ is the closest to } x_i \} \\ \text{size } n_l^t \\ y_l^{t+1} = \frac{1}{n_l^t} \sum_{i \in I_l^t} x_i \end{array} \right.$$

Output: $y_1^T \dots y_k^T$. \leftarrow Final centroids

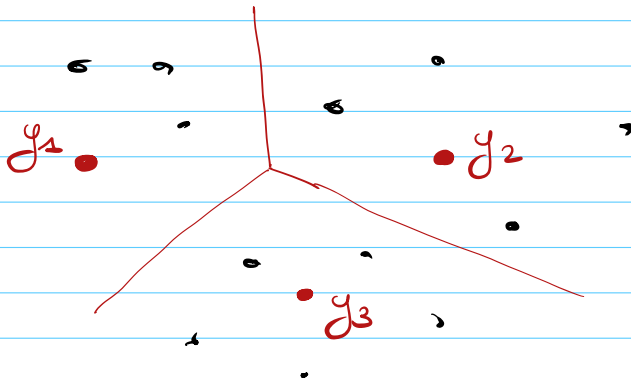


Prop: Lloyd's algorithm
 = Newton's method on $F_{k,x}$.

Recall: Newton's method $F: \mathbb{R}^d \rightarrow \mathbb{R}$

$$\rightarrow y' = y - (\nabla^2 F(y))^{-1} \nabla F(y) \parallel$$

Proof:



$$F_{k,x}(y_1, \dots, y_k) = \frac{1}{n} \sum_{i=1}^n \min_{j=1 \dots k} \|y_j - x_i\|^2$$

$$\nabla_{y_1} F_{k,x} = \begin{pmatrix} 0 \\ \frac{2}{n} \sum_{i \in I_1} x_i \\ 0 \end{pmatrix} = \frac{2}{n} \sum_{l=1}^k \sum_{i \in I_l} \|x_i - y_l\|^2$$

$$\nabla_{y_1} F_{k,x}(\dots) = \frac{1}{n} \sum_{i \in I_1} 2(y_1 - x_i)$$

$$= \frac{2}{n} (n_1 y_1 - \sum_{i \in I_1} x_i)$$

$$= \frac{2n_1}{n} \left(y_1 - \underbrace{\frac{1}{n_1} \sum_{i \in I_1} x_i}_{\tilde{y}_1} \right)$$

$$\Rightarrow \nabla F_{k,*}(-) = \begin{pmatrix} \nabla_{y_1} \\ \vdots \\ \nabla_{y_k} \end{pmatrix} = \frac{1}{2} \begin{pmatrix} n_1(y_2 - y_1) \\ \vdots \\ n_k(y_k - y_k) \end{pmatrix}$$

$$\Rightarrow \nabla^2 F_{k,*}(-) = \begin{pmatrix} \nabla_{y_1}^2 & & & \\ & \nabla_{y_1 y_2}^2 & & \\ & & \nabla_{y_2 y_2}^2 & \\ & & & \ddots \\ & & & & \nabla_{y_k}^2 \end{pmatrix}$$

$$= \frac{1}{2} \begin{pmatrix} n_1 I_d & 0 & & \\ 0 & n_2 I_d & & \\ 0 & & \ddots & \\ & & & n_k I_d \end{pmatrix}$$

$$\text{so } \nabla^2 F_{k,*}^{-1}(-) = \frac{1}{2} \begin{pmatrix} 1/n_1 & & 0 \\ & \ddots & \\ 0 & & 1/n_k \end{pmatrix}$$

$$(\nabla^2 F_{k,*}(-))^{-1} \nabla F_{k,*}(-) = \begin{pmatrix} y_1 - \bar{y}_1 \\ \vdots \\ y_k - \bar{y}_k \end{pmatrix}$$

\Rightarrow Newton's Method:

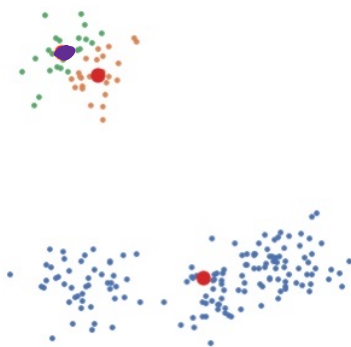
$$\begin{pmatrix} y_1 \\ \vdots \\ y_k \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_k \end{pmatrix} - \begin{pmatrix} y_1 - \bar{y}_1 \\ \vdots \\ y_k - \bar{y}_k \end{pmatrix} = \begin{pmatrix} \bar{y}_1 \\ \vdots \\ \bar{y}_k \end{pmatrix} = \text{One step of Lloyd's Algorithm.}$$



⇒ Newton's method on a
Non convex Function.

HW 8

↳ May not converge with initialization
too far away from optimum.

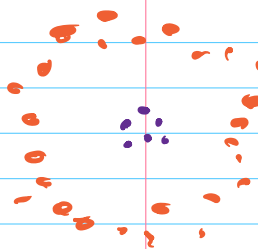


⇒ To Find a good initialization in
practice: k-means ++.

↳ sklearn default method.

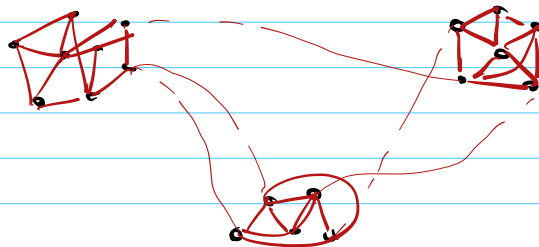
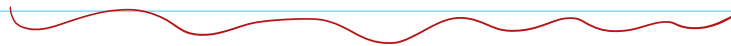
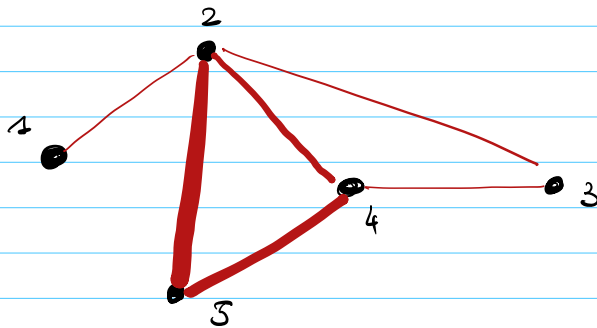
Summary: - most standard clustering algorithm

- will not work if clusters have "complex geometry"
- Restricted to $x_i \in \mathbb{R}^d$



② Spectral Clustering

→ input = weighted graph representing similarities between observations.

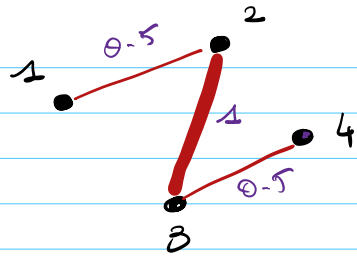


Three
Clusters.

Def: A weighted graph G with

- } n vertices
- } weights $W = (W_{ij})_{i,j}$ $n \times n$ symmetric matrix

$W_{ij} \geq 0$

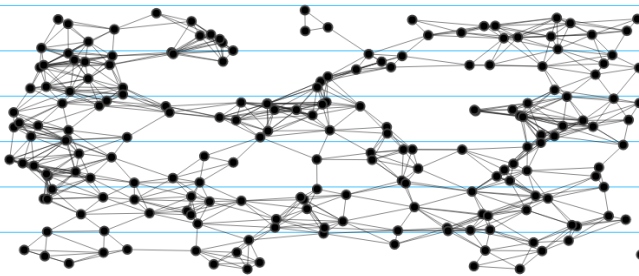


$$W = \begin{bmatrix} 1 & 0.5 & & \\ 0.5 & 1 & 1 & \\ & & 1 & 1 & 0.5 \\ & & & 0.5 & 1 \end{bmatrix}$$

Examples

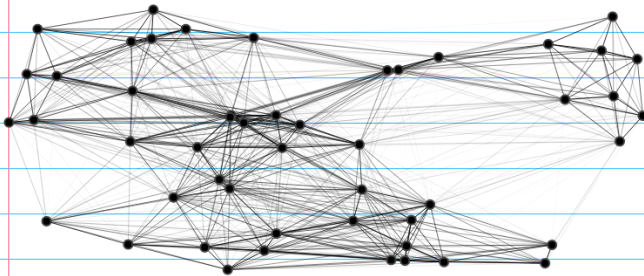
- (Non-weighted) graph:
$$\begin{cases} W_{ij} = 0 \text{ or } 1 \\ W_{ii} = 1 \end{cases}$$

- ϵ -Neighborhood graph



$$W_{ij} = \begin{cases} 1 & \text{if } \|x_i - x_j\| \leq \epsilon \\ 0 & \text{otherwise} \end{cases}$$

- Gaussian weights:



$$W_{ij} = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$$

Def: G weighted graph

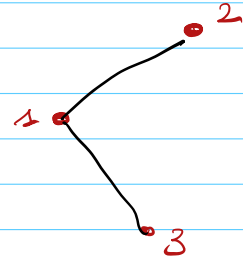
- **Neighbors** of i : j such that $w_{ij} > 0$

$\leadsto i \sim_{G} j$

- **Degree** D_i of i :

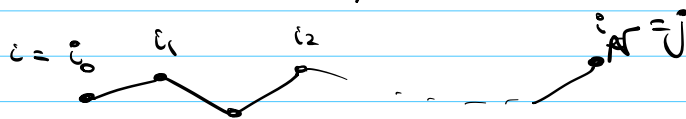
$$D_i = \sum_{j=1}^n w_{ij}$$

$$D = \begin{pmatrix} D_1 & & 0 \\ & \ddots & \\ 0 & & D_n \end{pmatrix} = \text{degree matrix}$$



- The vertices i and j are **connected**

if there is a path

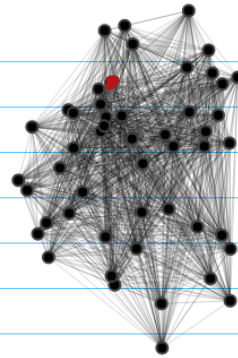
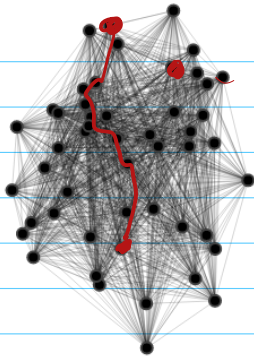


with $i_e \sim i_{e+1}$.

- A **connected component** is a set \mathcal{C} of vertices such that

$\forall i, j \in \mathcal{C}$ i and j connected

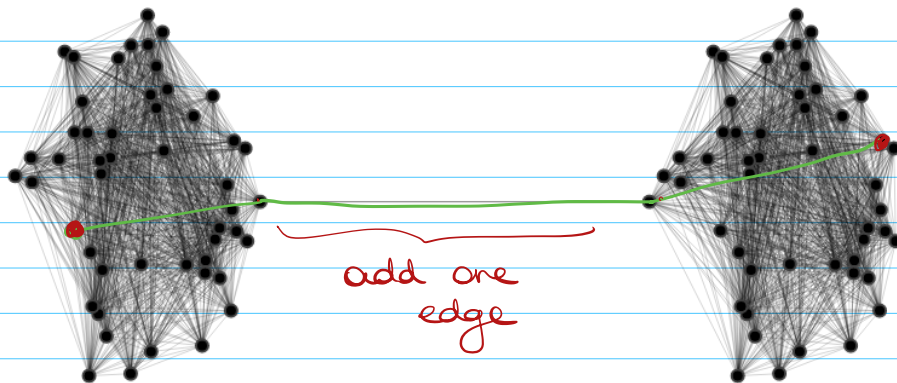
$\forall i \in \mathcal{C}, j \notin \mathcal{C}$ i and j NOT connected.



Two connected components

Idea 1: Clusters = connected components.

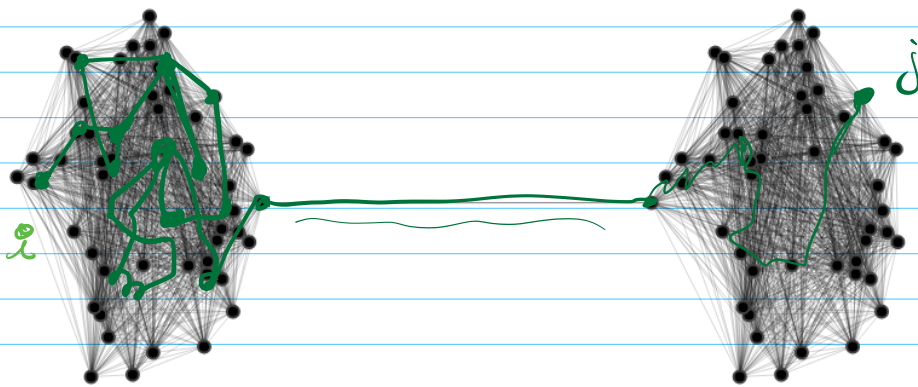
Problem: Not a robust notion



→ only one connected component.

Idea 2: Make the notion of being connected quantitative.

RANDOM WALK ON THE GRAPH.

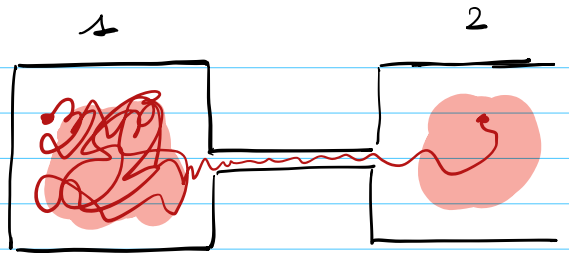


- ① Start at i_0 .
- ② Select a random neighbour i_1 of i_0 .
- ③ Repeat.
- ④ Stop when j is reached.

$$i = i_0 \sim i_1 \sim i_2 \dots \sim i_k = j$$

if k is large : i and j are
"almost" disconnected.

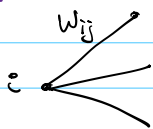
Analogy :



How much time does a molecule in a gas take to go from box 1 to box 2?

→ Probe to go from i to j in 1 step:

$$\sum_j Q_{ij} = \frac{1}{D_i} \sum_j W_{ij} = 1$$



$$Q_{ij} = \frac{1}{D_i} W_{ij}$$

↳ so that $\sum_j Q_{ij} = 1$

$Q = (Q_{ij})$ is the **probability transition matrix** of the walk.

$$\begin{pmatrix} 1-Q_{11} & -Q_{12} & -Q_{13} \\ & 1-Q_{22} & \\ & & \ddots \end{pmatrix} \begin{pmatrix} D_1 & & \\ & D_2 & \\ & & \ddots \end{pmatrix}$$

→ $L = I_n - Q =$ **Laplacian of \mathcal{G}** .

⇒ Spectral properties of L contain information about the geometry of \mathcal{G} .

$$A \subseteq \{1, \dots, n\} \quad e_A = (0 \dots 0 \overset{\text{elements of } A}{\underset{\uparrow}{1}} 0 \dots 0 \underset{\uparrow}{1} \dots) \in \mathbb{R}^n$$

$$(e_A)_i = \begin{cases} 1 & \text{if } i \in A \\ 0 & \text{otherwise.} \end{cases}$$

$$A = \{1\} \quad e_A = (1 \circ)$$

$$A = \{1, 2\} \quad e_A = (1 \ 1 \circ)$$

Prop:

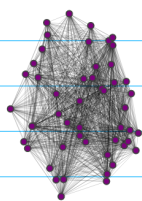
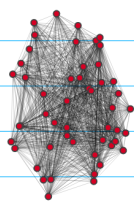
① L has nonnegative eigenvalues.

② 0 is an eigenvalue of L .

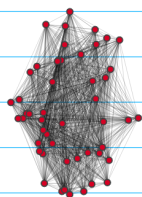
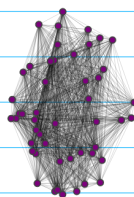
③ The multiplicity of 0 is the number k of connected components of \mathcal{G} .

A basis of the eigenspace is


$$\{e_{c_1}, \dots, e_{c_k}\} \quad \text{connected component}$$



$$e_{c_1} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 0 \end{pmatrix}$$



$$e_{c_2} = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 1 \\ 1 \end{pmatrix}$$

proof:  L is NOT symmetric!

^{symmetric}
Laplacian $\tilde{L}' = D^{1/2} L D^{-1/2} = I_n - D^{1/2} D^{-1} W D^{-1/2}$
 $= I_n - D^{-1/2} W D^{-1/2} \leftarrow$ is symmetric

$$L u = d u \iff L' v = d v \quad v = D^{1/2} u$$

\Rightarrow L and L' have the same eigenvalues.

Let $v \in \mathbb{R}^n$

$$v^T L' v =$$

$$= \frac{1}{2} \sum_{1 \leq i, j \leq n} w_{ij} \left(\frac{v_i}{\sqrt{d_i}} - \frac{v_j}{\sqrt{d_j}} \right)^2 \geq 0.$$

② Take $v = (\sqrt{d_1}, \dots, \sqrt{d_n}) \Rightarrow v^T L' v = 0$
 $\Rightarrow 0$ eigenvalue of L'
 $\Rightarrow 0$ eigenvalue of L.

③

v eigenvector of L'

$$[v = D^{1/2}u]$$

$\Leftrightarrow u$ eigenvector of L .

\Rightarrow When do we have

$$0 = v^T L' v = \frac{1}{2} \sum_{1 \leq i, j \leq n} w_{ij} \left(\frac{v_i}{\sqrt{d_i}} - \frac{v_j}{\sqrt{d_j}} \right)^2$$

$$= \frac{1}{2} \sum_{1 \leq i, j \leq n} w_{ij} (u_i - u_j)^2 \Rightarrow \begin{cases} u_i = u_j \\ \text{if } w_{ij} > 0 \end{cases}$$

$k=1$: one connected component.

• $u = (1 \dots 1)$ is an eigenvector.

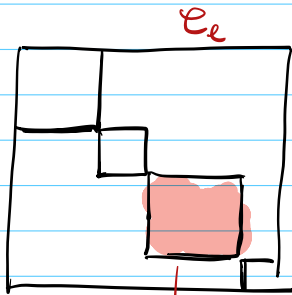
• if $u \in \mathbb{R}^n$ s.t. $v^T L' v = 0$

Path $i = b \sim i_1 \dots \sim i_k = j$

$$u_i = u_{i_1} = u_{i_2} \dots = u_{i_k} = u_j$$

$\Rightarrow u = (c, c, \dots, c)$

$k > 1$:



subgraph that
is connected

$\Rightarrow eee$ is an
eigenvector.



Summary:

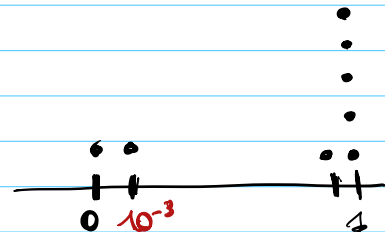
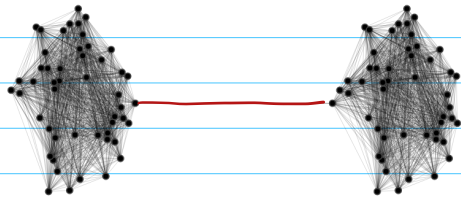
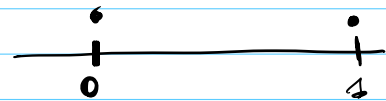
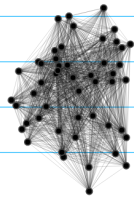
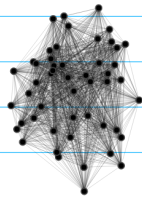
Multiplicity of 0

↕
Number of
connected
components

Associated Eigenvectors

↕
Connected components

What if we add an edge?



↪ Stability of the spectrum with respect to small perturbations.

k clusters = k small eigenvalues

↳ associated eigenvectors give the clusters.

SPECTRAL CLUSTERING

- 1- Compute L
- 2- Compute the k first eigenvalues / eigenvectors of L .
 $0 = d_1 \leq d_2 \leq \dots \leq d_k$
 $v_1 \quad v_2 \quad \dots \quad v_k$ eigenvectors. ($\in \mathbb{R}^n$)
- 3- Let $x_i = (\langle v_1, e_i \rangle, \dots, \langle v_k, e_i \rangle) \in \mathbb{R}^k$
 $= (0, 1, 0)$
 $=$ represents i .
- 4- Apply k -means on (x_1, \dots, x_n) .

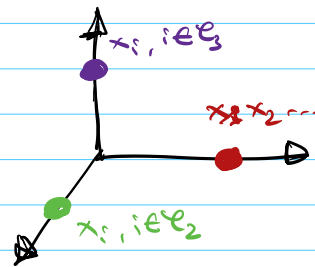
\Rightarrow if k connected components:

$$v_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad v_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad v_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

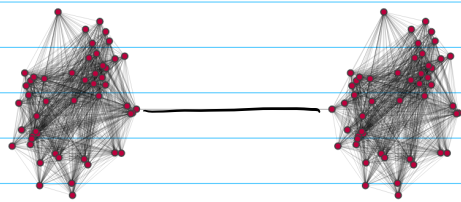
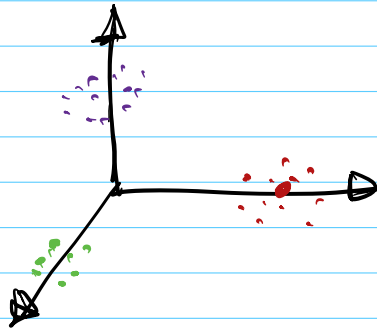
$$i \in \mathcal{C}_1 \quad x_i = (1, 0, 0)$$

$$i \in \mathcal{C}_2 \quad x_i = (0, 1, 0)$$

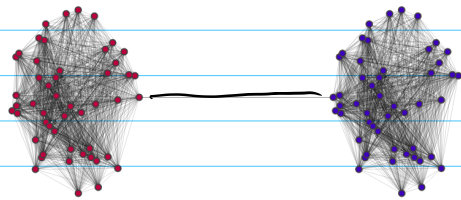
$$i \in \mathcal{C}_3 \quad x_i = (0, 0, 1)$$



→ With [~]approximate[~] connected components:
clusters



$$d = 0$$



$$d = 10^{-3}$$