

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.



Q: How should the centroids be defined?

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

The best way to summarize \mathcal{X} using one point:

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

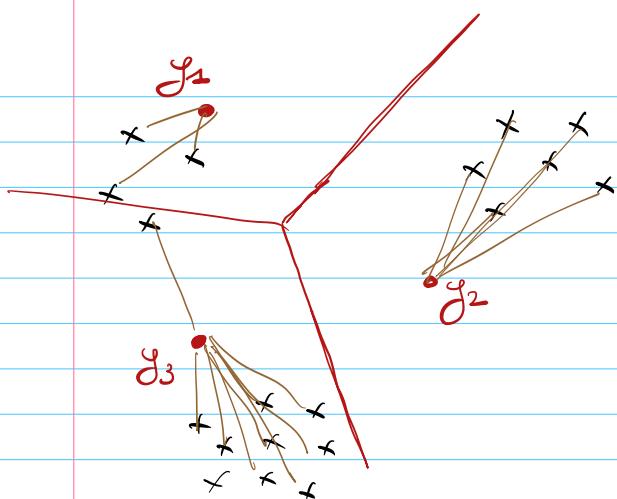
Actually, y^* is the minimizer of

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

proof:

Q: For $k > 1$: Generalization

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

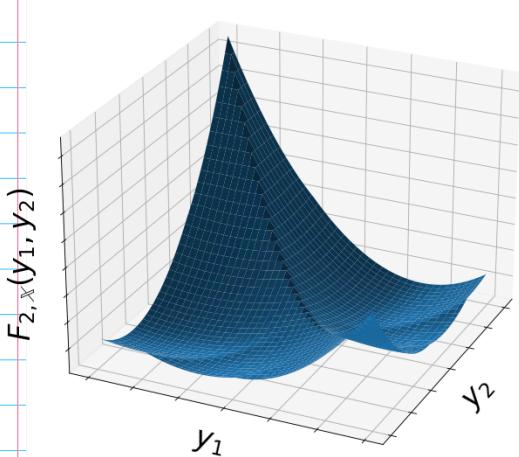


Choose (g_1, g_2, g_3)

to minimize

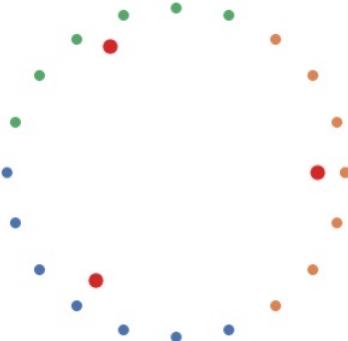
$$\sum (-)^2$$

\Rightarrow The k -means of X is the minimizer (g_1^+, \dots, g_k^+) of $F_{k,X}$.
 ↳ also called centroids



The function $F_{k,X}$

is Not convex.
 $(k > 2)$



Several k -means
may exist.

How to compute the k -means?

Lloyd's ALGORITHM

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

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

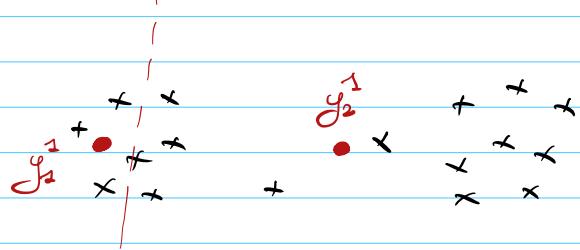
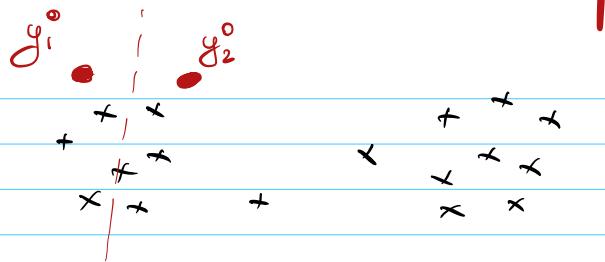
$\forall l=1 \dots k$:

$$I_l^t = \{i : y_l^t \text{ is the closest to } x_i\}$$

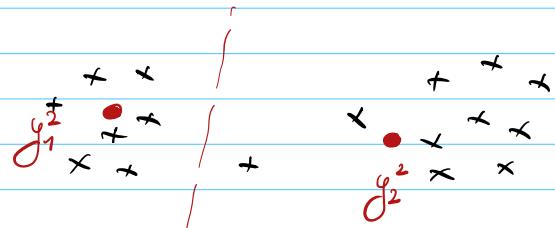
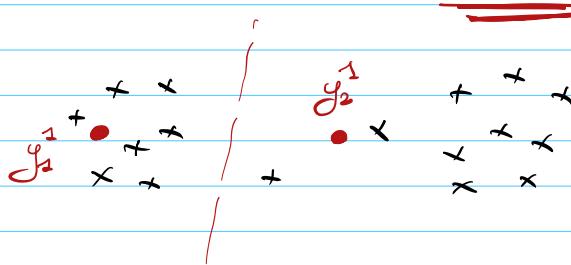
size n_l^t

$$y_l^{t+1} = \frac{1}{n_l^t} \sum_{i \in I_l^t} x_i$$

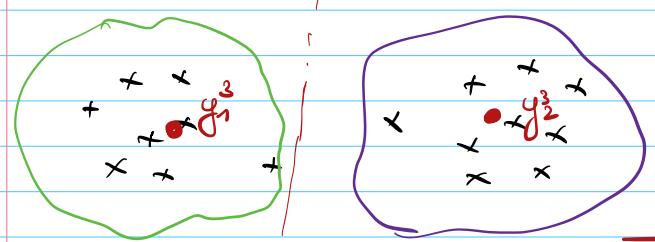
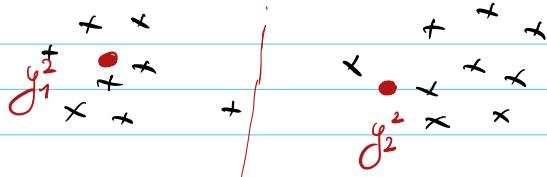
Output: $y_1^T \dots y_k^T$. \leftarrow final centroids



Step 1



Step 2



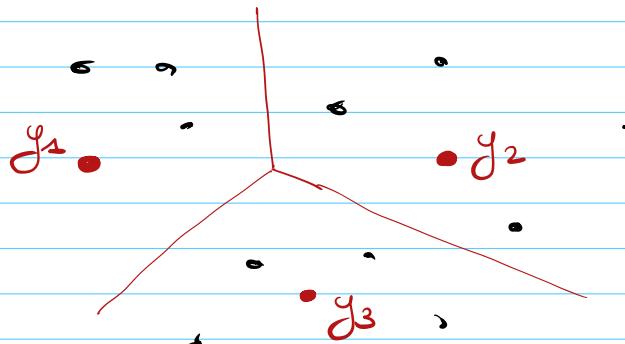
END

$\underbrace{\text{Prop}}$: Lloyd's algorithm
 = Newton's method on $F_{k,\infty}$.

Recall: Newton's method

$$x' = x - (\nabla^2 F(x))^{-1} \nabla F(x)$$

proof:



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

$$= \frac{1}{n} \sum_{\ell=1}^k \sum_{i \in I_\ell} \|x_i - y_\ell\|^2$$

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

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

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

$$\Rightarrow \nabla \tilde{F}_{k,*}(-) = \begin{pmatrix} \nabla_{y_1} \\ \vdots \\ \nabla_{y_k} \end{pmatrix} = \frac{2}{n} \begin{pmatrix} n_1(y_1 - \tilde{y}_1) \\ \vdots \\ n_k(y_k - \tilde{y}_k) \end{pmatrix}$$

$$\Rightarrow \nabla^2 \tilde{F}_{k,*}(-) = \begin{pmatrix} \nabla^2_{yy_1} & & \\ & \nabla^2_{y_1 y_2} & \\ & & \ddots \\ & & & \nabla^2_{y_k y_k} \end{pmatrix}$$

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

$$\text{and } \nabla^2 \tilde{F}_{k,*}^{-1}(-) = \frac{n}{2} \begin{pmatrix} n_1 & 0 \\ 0 & \ddots \\ & & n_k \end{pmatrix}$$

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

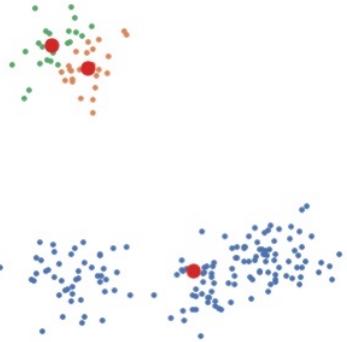
\Rightarrow Newton's Method:

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

\Rightarrow Newton's method or a
Non convex Function.

HW 8

\hookrightarrow May not converge with initialization
too far away from optimum.

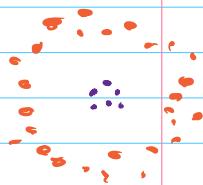


\Rightarrow To Find a good initialization in
practice: k-means ++.

\hookrightarrow sklearn default method.

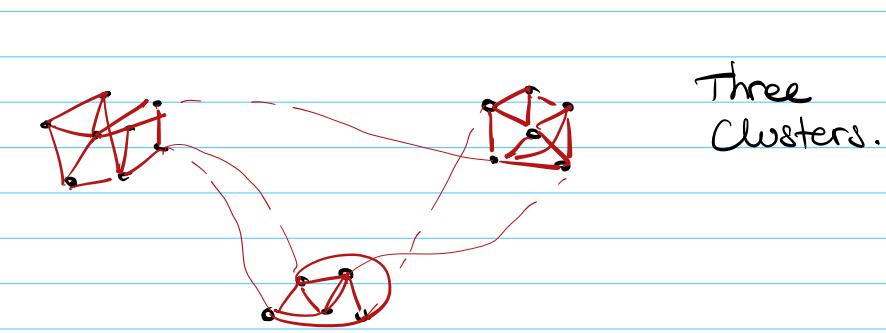
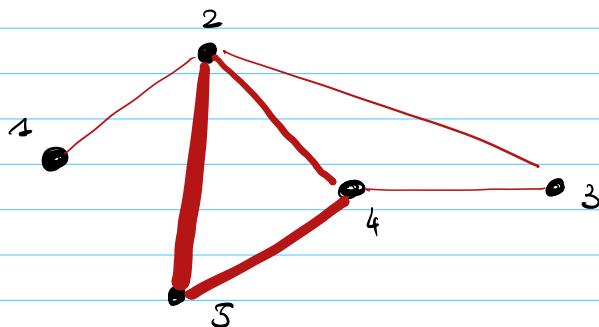
Summary: - most standard clustering algorithm

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



② Spectral Clustering

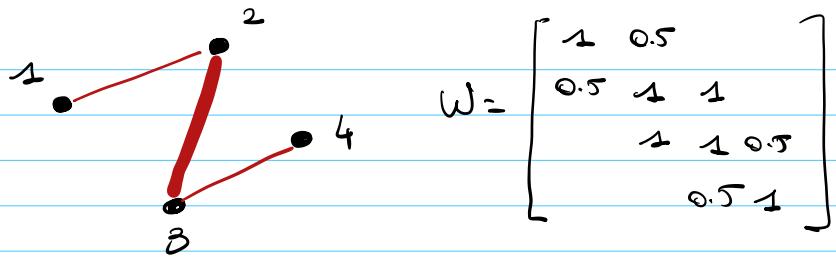
→ input = weighted graph representing similarities between observations.



Three Clusters.

Def: A weighted graph G with

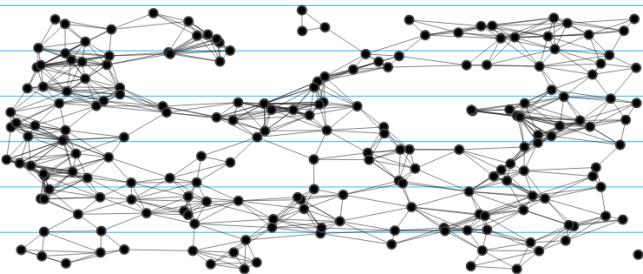
$\left\{ \begin{array}{l} n \text{ vertices} \\ \text{weights } W = (W_{ij})_{ij} \text{ } n \times n \text{ symmetric} \\ \text{matrix} \\ W_{ij} \geq 0 \end{array} \right.$



Examples

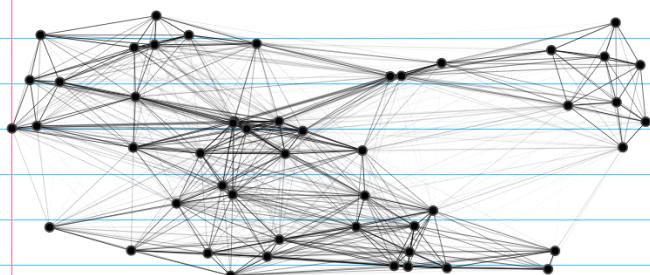
- (Non-weighted) graph: $\begin{cases} W_{ij} = 0 \text{ or } 1 \\ W_{ii} = 1 \end{cases}$

- ϵ -Neighborhood graph



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

- Gaussian weights:



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

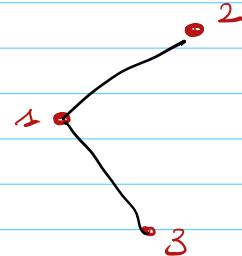
Def: If weighted graph

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

$\Rightarrow i \sim j$.

- Degree D_i of i :

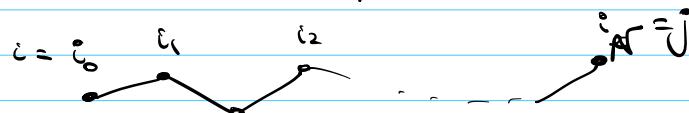
$$D_i = \sum_{j=1}^n w_{ij}$$



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

- The vertices i and j are connected

if there is a path

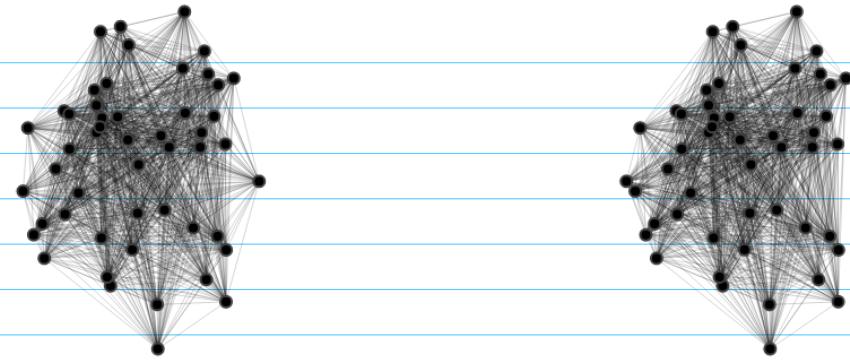


with $i_e \sim i_{e+1}$.

- A connected component is a set C of vertices such that

$\forall i, j \in C$ i and j connected

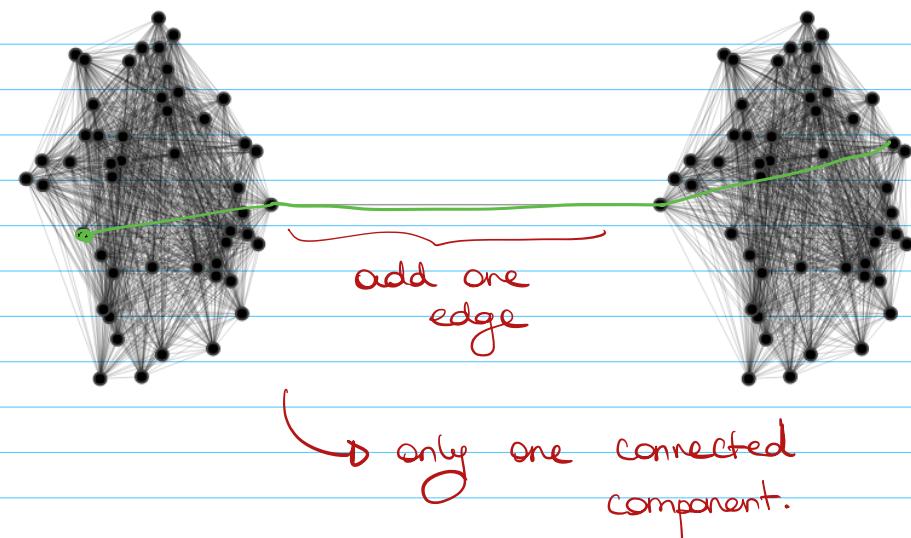
$\forall i \in C, j \notin C$ i and j NOT connected.



Two connected
components

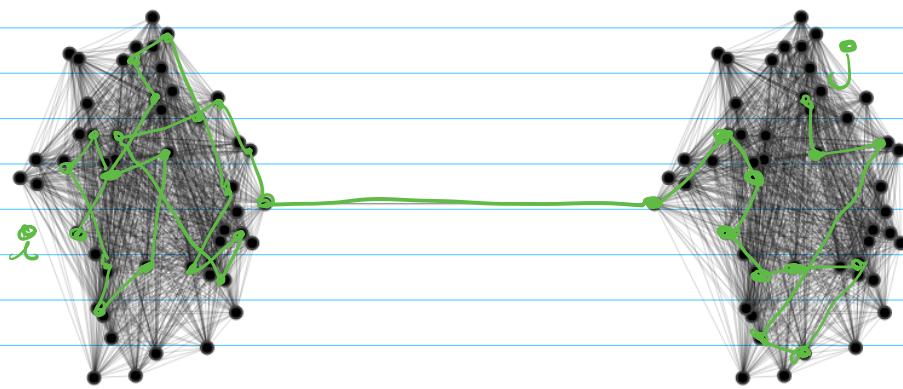
Idea 1 : Clusters = connected components.

Problem: Not a robust notion



Idea 2 : Make the notion of being connected quantitative.

RANDOM WALK ON THE GRAPH.

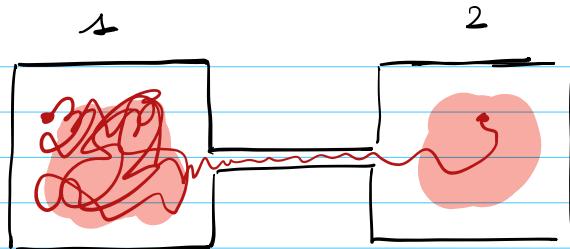


- ① Start at i_0 .
- ② Select a random neighbour i_1 of i_0 .
- ③ Repeat.
- ④ Stop when j is reached.

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

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:

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

$\therefore \text{so that } \sum_j Q_{ij} = 1$

$Q = (Q_{ij})$ is the probability transition matrix

of the walk.

→ $L = I_n - Q$ = Laplacian of G .

⇒ Spectral properties of L contain
information about the geometry of G .

$$A \subseteq \{1, \dots, n\} \quad e_A = (0, \dots, 0, \overset{\text{elements of } A}{1}, 0, \dots, 0)$$

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

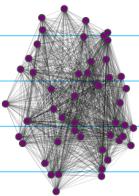
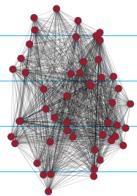
Prop:

- (1) L is symmetric positive semi-definite.
- (2) 0 is an eigenvalue of L .
- (3) The multiplicity of 0 is the number k of connected components of G .

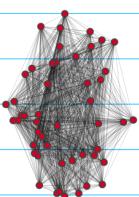
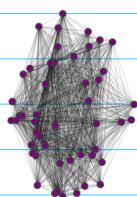
A basis of the eigenspace is

$$\{e_{C_1}, \dots, e_{C_k}\}$$

\hookrightarrow connected component



$$e_{C_1} = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$



$$e_{C_2} = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

proof:



L is NOT symmetric!

\sim symmetric
Laplacian $(L' = D^{1/2} L D^{-1/2} = I_n - D^{1/2} W D^{-1/2})$
 $= I_n - D^{1/2} W D^{-1/2}$ ← is symmetric

$$Lu = dv \quad \Leftrightarrow \quad L'v = dv \quad v = D^{1/2}u$$

⇒ 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.

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

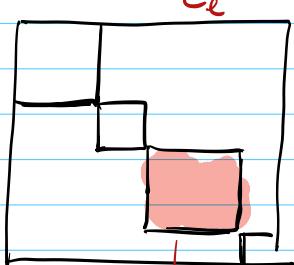
• if $u \in \mathbb{R}^n$ s.t. $v^T L' v = 0$

Path $i_0 \sim i_1 \sim \dots \sim i_k = j$

$$u_{i_0} = u_{i_1} = u_{i_2} = \dots = u_{i_k} = u_j$$

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

$k \geq 1$:



subgraph that
is connected

e_{ee} is an
eigenvector.



Summary:

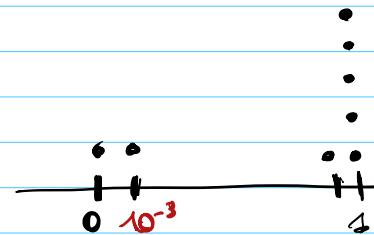
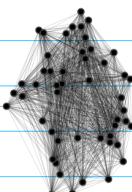
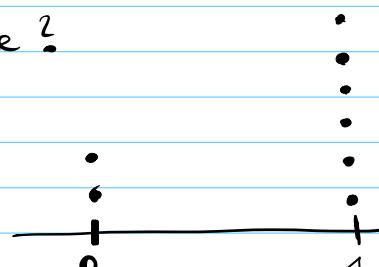
Multiplicity of 0

↔
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 .

$$0 = \lambda_1 \leq \lambda_2 \dots \leq \lambda_k$$

v_1, v_2, \dots, v_k eigenvectors. ($\in \mathbb{R}^n$)

3- Let $x_i = (\underbrace{\langle v_1, e_i \rangle, \dots, \langle v_k, e_i \rangle}_{= (0, 1, 0)}, \dots)$ represents i .

4- Apply k -means on (x_1, \dots, x_n) .

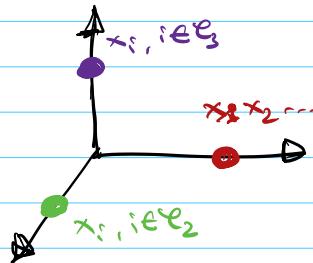
→ if k connected components:

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

$$i \in C_1 \quad x_i = (1, 0, 0)$$

$$i \in C_2 \quad x_i = (0, 1, 0)$$

$$i \in C_3 \quad x_i = (0, 0, 1)$$



With "approximate" connected components:

clusters

