

# Modern Clustering algorithms

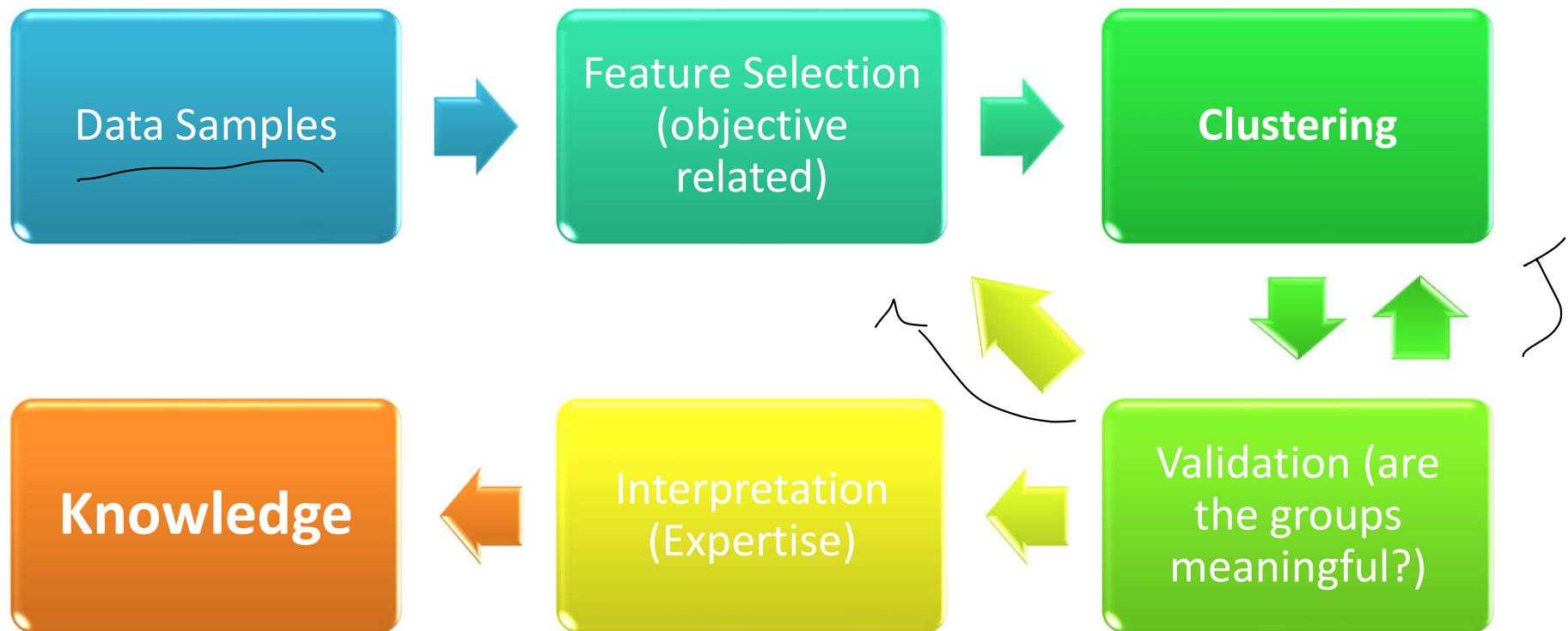
Unsupervised Machine Learning

# Clustering procedure

- Cardinality,
- Dimension, type

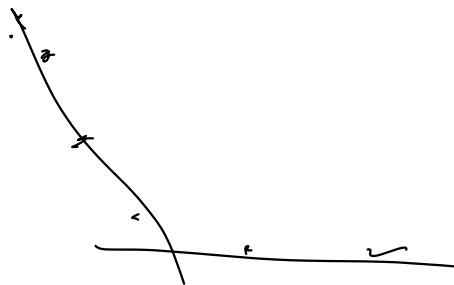
Feat. ranking,  
Dim. Red., Dist.

K-means,  
fuzzy c-means,  
hierarchical



# Some problems with the discussed algorithms

- K-means and fuzzy c-means are not well suited to deal with arbitrary shape (non-convex) clusters.
- Which intercluster distance shall we use for hierarchical clustering? (and can I choose?)
- What happens with noisy data sets?
- Which should be the final number of clusters?



→ what happens if  
"no elbow" in some  
plots?

# Clustering procedure

Cardinality,  
Dimension, type



Feat. ranking,  
Dim. Red., Dist.



*relatively modern*

“Modern”  
Clustering

Knowledge

Interpretation  
(Expertise)

Validation (are  
the groups  
meaningful?)



# Reviews on clustering algorithms

- There are hundreds of clustering algorithms, each of them with their advantages and withdrawals.
  - Jain, A. K., Murty, M. N., & Flynn, P. J. (1999). Data clustering: a review. *ACM computing surveys (CSUR)*, 31(3), 264-323.
  - Xu, R., & Wunsch, D. (2005). Survey of clustering algorithms. *Neural Networks, IEEE Transactions on*, 16(3), 645-678.
  - Jain, A. K. (2010). Data clustering: 50 years beyond K-means. *Pattern recognition letters*, 31(8), 651-666.
  - Xu, D., & Tian, Y. (2015). A Comprehensive Survey of Clustering Algorithms. *Annals of Data Science*, 2(2), 165-193.
  - Amit Saxena, et al (2017). A review of clustering techniques and developments, *Neurocomputing*, 267, 664-681.

# A selected list of clustering algorithms

scikit-learn

- K-means
- Fuzzy c-means
- Agglomerative hierarchical
- Spectral Clustering
- Gaussian Mixture
- Affinity Propagation
- Density Based Algorithms:
  - DBSCAN
  - Density Peaks
  - Mean-Shift clustering

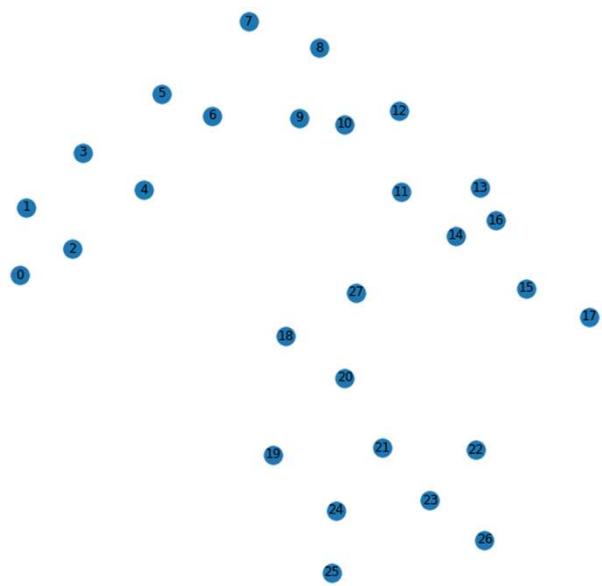
# A selected list of clustering algorithms

- K-means Classical
- Fuzzy c-means
- Agglomerative hierarchical
- Spectral Clustering Modern
- Gaussian Mixture
- Affinity Propagation
- Density Based Algorithms:
  - DBSCAN
  - Density Peaks
  - Mean-Shift clustering

# Spectral clustering

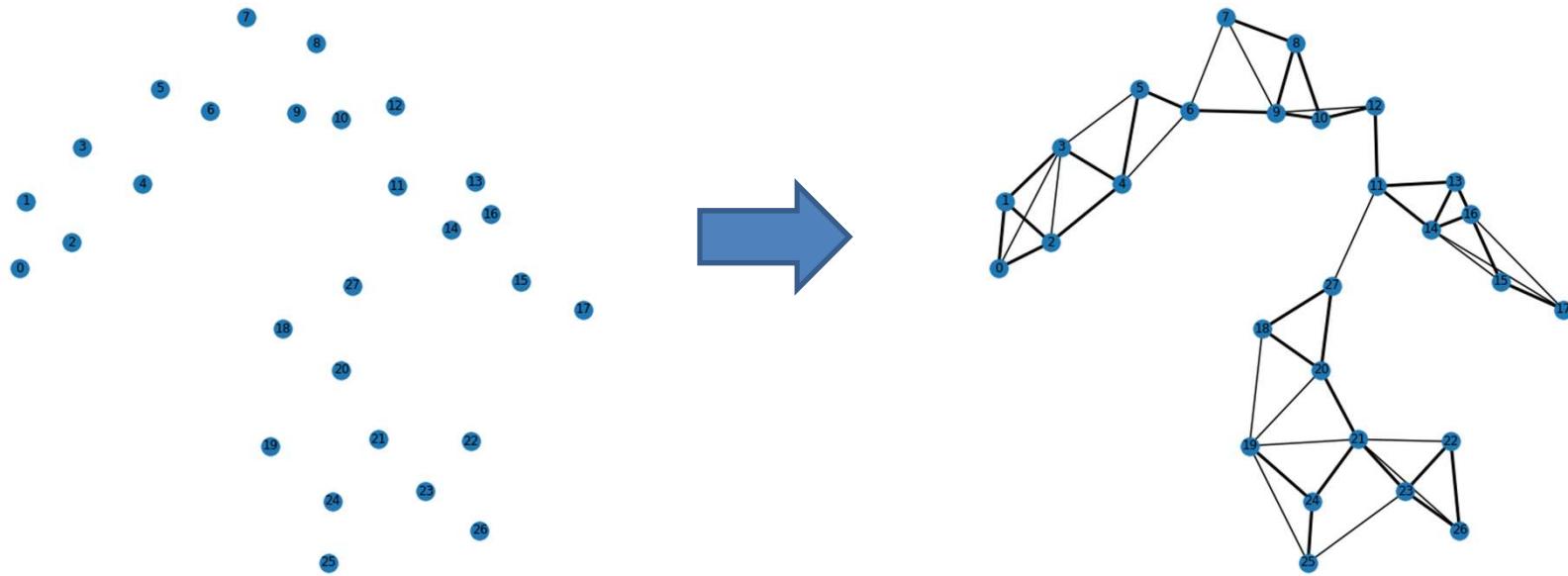
Modern clustering algorithms (I)

# A different paradigm



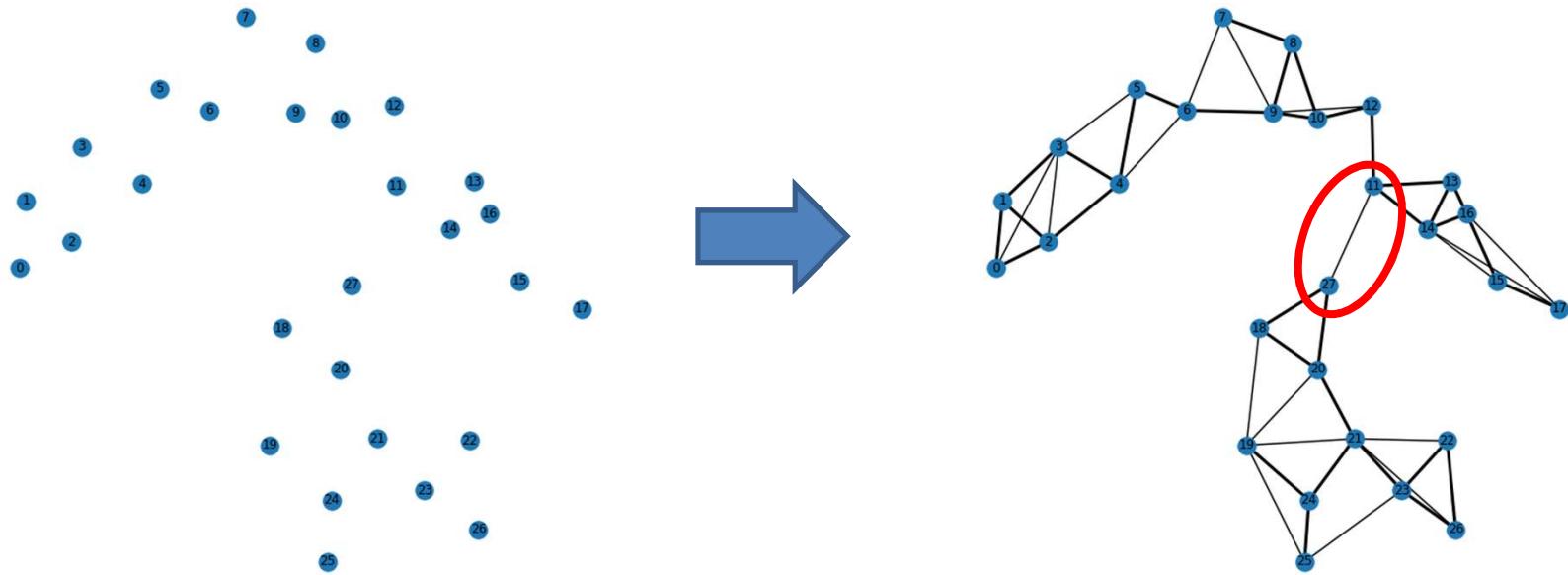
instead of directly applying  
algorithms on the  
data, we first transform  
it into a graph

# A different paradigm



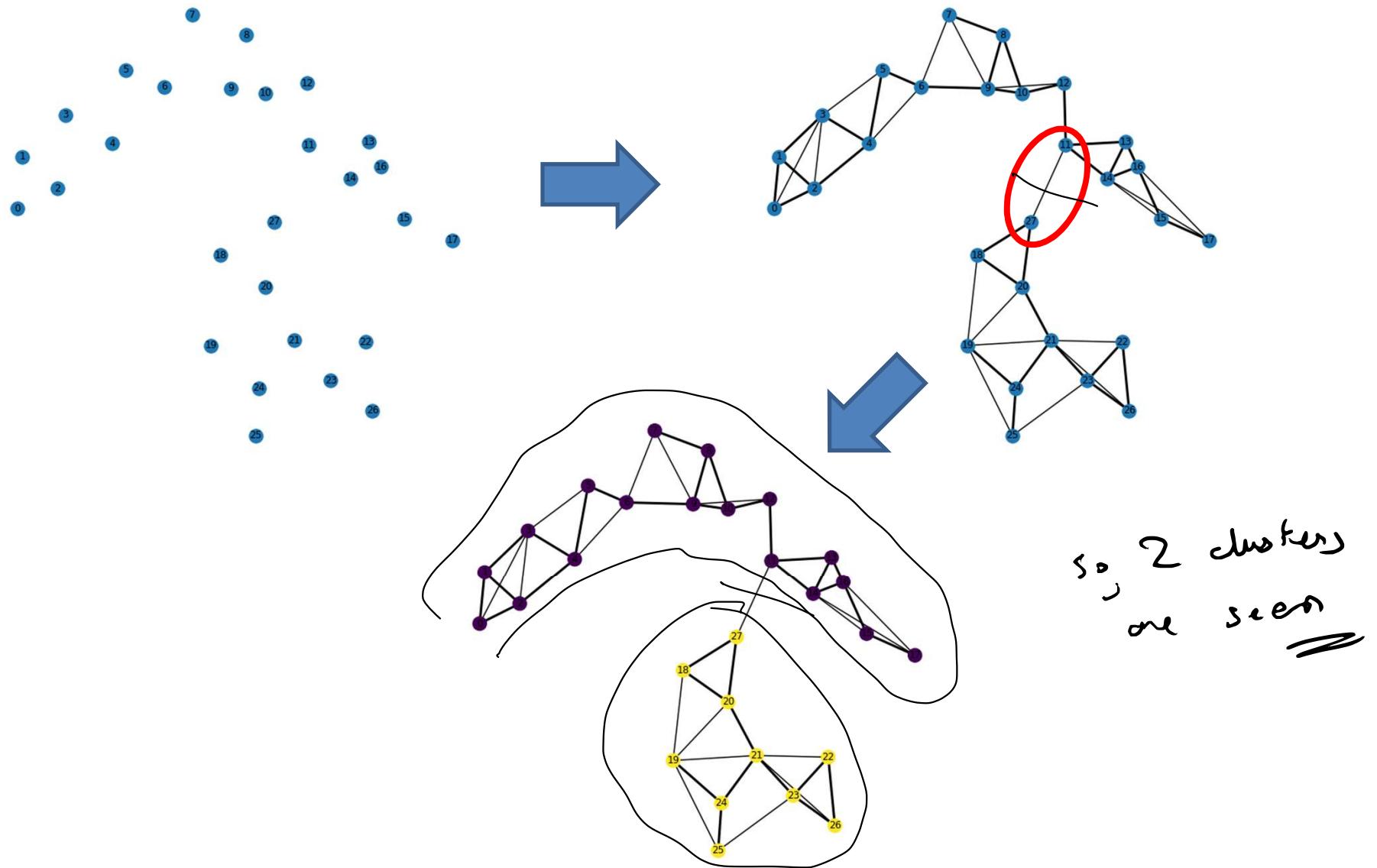
① Transform into a graph

# A different paradigm

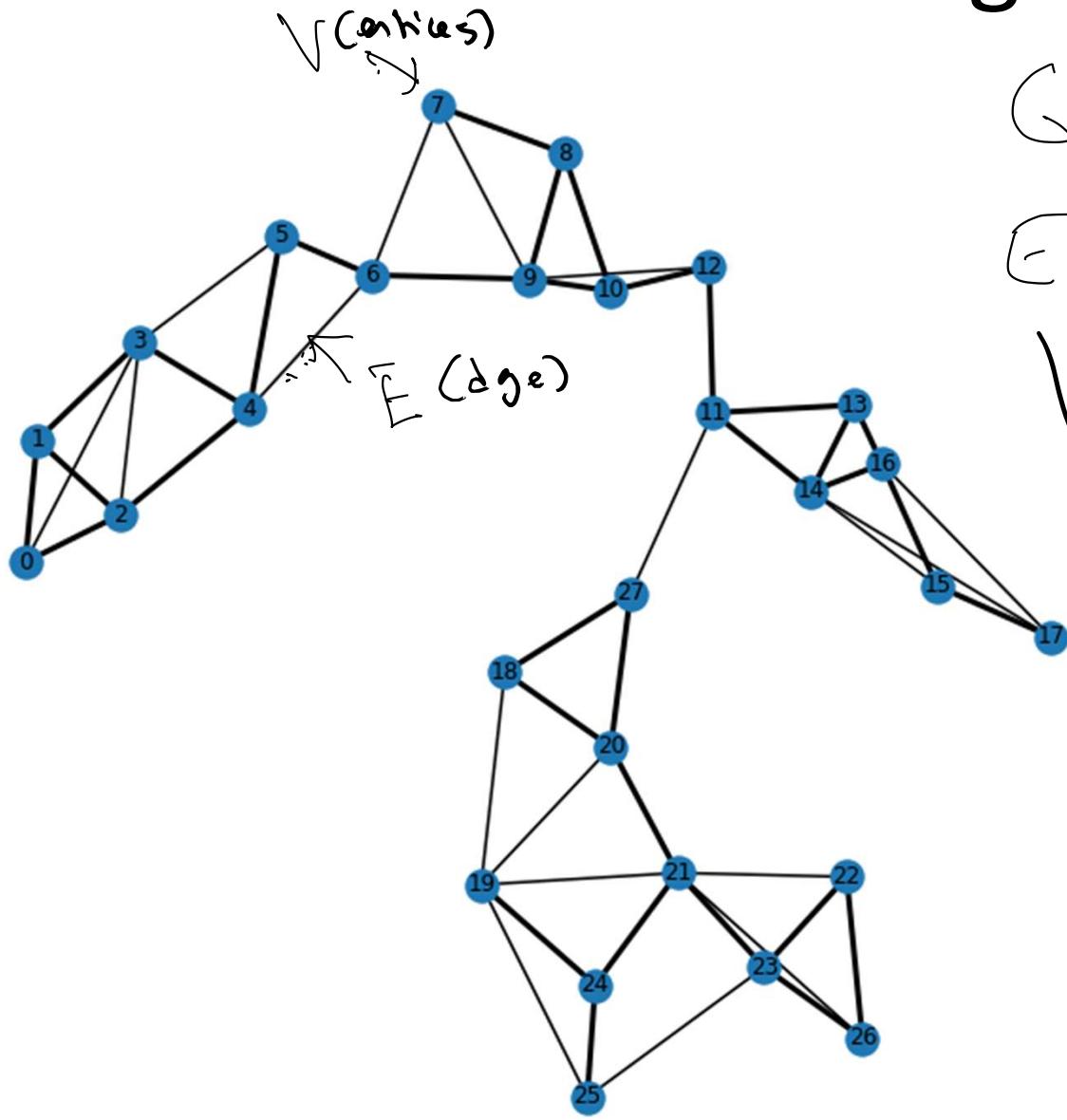


② Optimal way of ~~cutter~~  
cutting the graph

# A different paradigm



# Undirected Weighted Graph

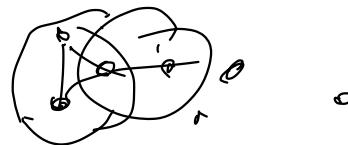


$G(E, V) \rightarrow$  graph defined by edges & vertices  
 $E \{ (i, j), S_{i,j} \geq 0 \}$   
→ no arrows  
→ no weights for edges

each edge is defined by the 2 vertices ( $i$  &  $j$ )  
 $t = \text{weight}$ .  
if weight = 0 i.e. if one is not connected.

# Obtaining the (similarity) graph

- **$\epsilon$ -ball graph:** Join vertices within a radius  $\epsilon$ .



- **$k$ -NN graph:** Join the  $k$ -Nearest Neighbor vertices.

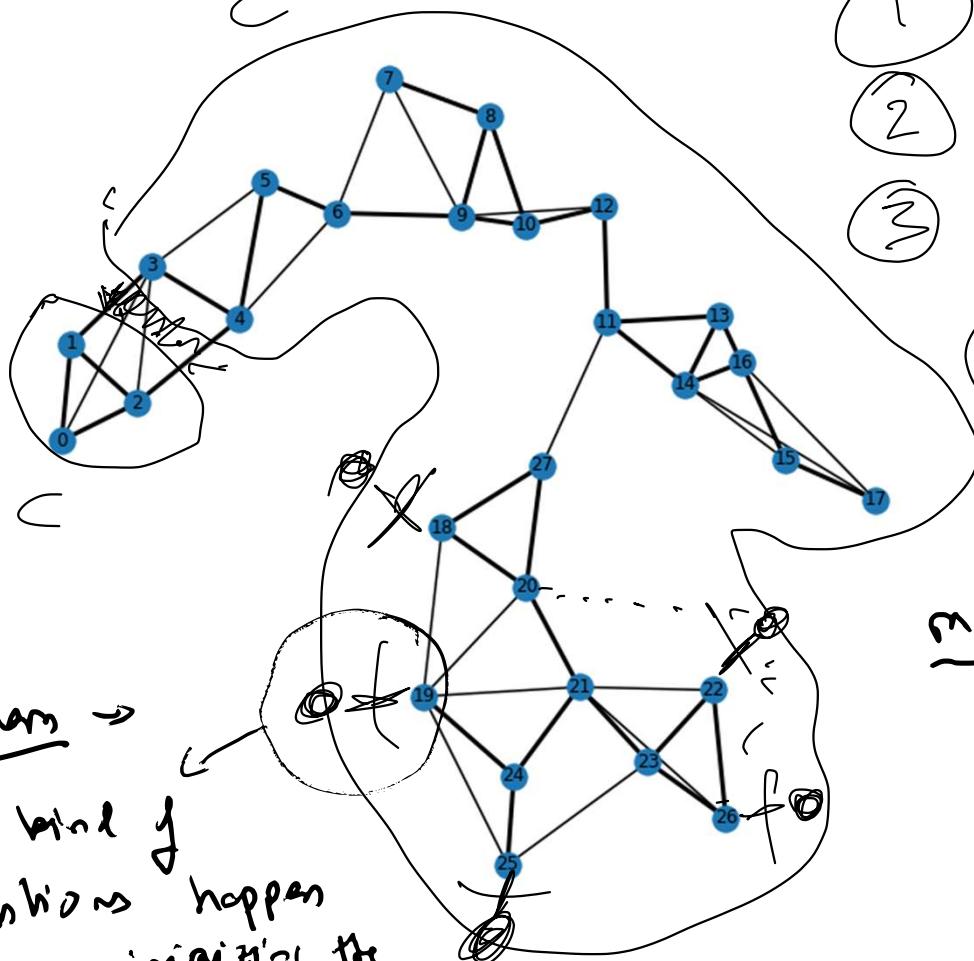


- **Fully connected graphs:** Join all the vertices,

for instance:  $S_{ij} = e^{\frac{-d_{ij}^2}{2\sigma^2}}$

---

# The k-way cut: a naive approach



problem →  
these kind of positions happen  
when minimizing the  
cut(5) & we get  
single points / noise in the result

- 1 Cluster  $C$  set of vertices
- 2  $\bar{C}$  all the vertices  $\notin C$
- 3  $\text{Cut}(C, \bar{C}) = \sum_{i \in C} \sum_{j \in \bar{C}} S_{i,j}$
- 4  $\text{Cut}(C_1, C_2, \dots, C_k) = \frac{1}{2} \sum_e \text{Cut}(C_e, \bar{C}_e)$
- 5 Minimize k-way Cut

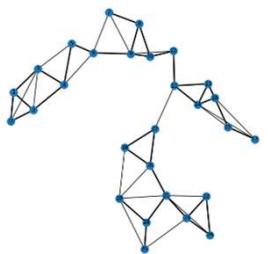
more explanation of above approach :-

- 1 → define a cluster
- 2 → opposite of ' $C$ ' - all vertices not belonging to  $C$
- 3 → define the cut  $\bar{C}$  & all vertices belonging to  $C \neq \bar{C}$  (complement of  $C$ ) as a sum.
- 4 → generalize (3) to ' $k$ ' clusters

$\rightarrow$   $S_0$  to  $n_{\text{cupp}}$ : the problem is we don't take size of clusters into account.

## Balancing the k-way cut: Ratio Cut

- $|C_\ell|$  = Number of vertices of  $C_\ell$
- Ratio Cut  $\underline{\ell C_1, C_2, \dots, C_K} = \sum_\ell \frac{\text{Cut}(C_\ell, \bar{C}_\ell)}{|C_\ell|}$ 
  - define a different "cut" → this designs less the clusters with larger population.



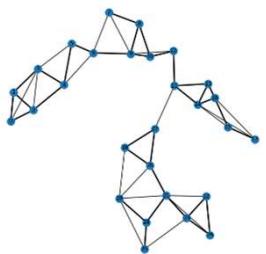
# Another way of taking into account the size of the clusters: Normalized cut

$$\rightarrow \text{Degree of vertex} := d_i = \sum_j S_{ij}$$

$$\rightarrow \text{Vol}(C_\ell) = \sum_{i \in C_\ell} d_i \quad (\text{Vol } \underset{\ell}{\underbrace{\text{of}}} \text{ a cluster})$$

$$\rightarrow N\text{cut}(C_1, C_2, \dots, C_k) = \sum_{\ell} \frac{\text{Cut}(C_\ell, \bar{C}_\ell)}{\text{Vol}(C_\ell)}$$

$\curvearrowright$  (Normalized cut)



→ Some def<sup>ns</sup> needed for next part :-

## The path to Spectral clustering: Matrices involved

(Similarity matrix)

$\$$  = Symmetric

$$n \times n$$
$$n \equiv \# \text{ of data points}$$

diagonal = 0,  
elements are  
 $s_{ij}$ .  
i.e. group  
is individ  
- ualized  
so it is  
symmetric

(Degree matrix)

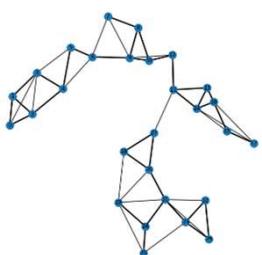
$$\Pi : \{ D_{ii} = d_i, D_{ij} = \frac{1}{\sqrt{d_j}} \text{ if } i \neq j \}$$

(Laplacian matrix)

$$L = \Pi - \$$$

diagonal terms  
= degree of matrix  
non-diagonal = 0

at the core  
of spectral clustering



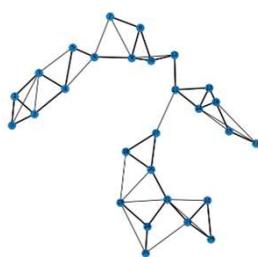
# The path to Spectral clustering: Important properties of the Laplacian Matrix

① Symmetric  $\mathbb{L} = \mathbb{D} - \mathbb{S}$ ,  $\therefore \mathbb{L}$  is sym.  
 $\begin{matrix} \uparrow \\ \mathbb{D}_{\text{sym}} \end{matrix} \quad \begin{matrix} \uparrow \\ \mathbb{S}_{\text{sym}} \end{matrix}$

② Positive semi-definite:  $\lambda_1 \geq \lambda_2 \geq \dots$   
this define it.  $v \in \mathbb{R}^n$  no compute  $v^T \mathbb{L} v$

③  $v$  as vector dim =  $n$

$$\begin{aligned}
 d_i &= \sum_j S_{ij} \\
 d_j &= \sum_i S_{ij} \\
 &= \frac{1}{2} \left( \sum_i v_i^2 d_i - 2 \sum_{ij} v_i v_j S_{ij} + \sum_j v_j^2 d_j \right) = \\
 &\quad \frac{1}{2} \sum_{ij} v_i^2 S_{ij} - 2 \sum_{ij} v_i v_j S_{ij} + \sum_{ij} v_j^2 S_{ij} = \\
 &\quad \frac{1}{2} \sum_{ij} S_{ij} (v_i^2 - 2 v_i v_j + v_j^2) = \frac{1}{2} \sum_{ij} S_{ij} (v_i - v_j)^2
 \end{aligned}$$

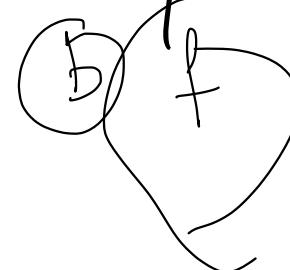


"indicates vector"  
 tells us if "single vertex belongs to  
 a cluster ( $C$ ) or its complement ( $\bar{C}$ )  
 let's define a vector  $f$   
 with following properties

## The RatioCut for $k=2$

(a) Minimize

$$\frac{\text{Cut}(C, \bar{C})}{|C|}$$



$$f_i = \sqrt{\frac{|\bar{C}|}{|C|}} \text{ if } i \in C$$

$$f_i = \sqrt{\frac{|C|}{|\bar{C}|}} \text{ if } i \in \bar{C}$$

$$(b) f^T \parallel f = \frac{1}{2} \sum_{ij} S_{ij} (f_i - f_j)^2 = \frac{1}{2} \sum_{\substack{i \in C \\ j \in \bar{C}}} S_{ij} \left( \sqrt{\frac{|\bar{C}|}{|C|}} - \sqrt{\frac{|\bar{C}|}{|C|}} \right)^2$$

$\approx (i) \& (j) \text{ belong to different clusters}$

$$+ \frac{1}{2} \sum_{\substack{i \in C \\ j \in C}} S_{ij} \left( -\sqrt{\frac{|C|}{|\bar{C}|}} - \sqrt{\frac{|C|}{|\bar{C}|}} \right)^2 =$$

$$\text{cut}(C, \bar{C}) \cdot \left( \frac{|\bar{C}|}{|C|} + \frac{|C|}{|\bar{C}|} + 2 \right) = n \text{ RatioCut}(C, \bar{C})$$

$C$  can be derived from some algebra

when  $\{i \in C, j \in C\} \neq \{i \in \bar{C}, j \in \bar{C}\}$ , the terms cancel out,

# The RatioCut for k=2

minimize

$$\frac{\text{Cut}(C, \bar{C})}{|C|}$$

equivalent

minimize

$$f^T K f \quad \text{sign: } \|f\| = \sqrt{n} \times f + 1$$

condition ①

condition ②  
↓

This avoids  
the trivial  
 $\text{sol}^2$ , where all  
values of  $f$   
are the same.

# The RatioCut for arbitrary k

→ K indicator vectors,  $\mathbb{H}$   $i = 1 \dots n$  (define "h")  
 ①  $H_{ij} = \frac{1}{\sqrt{|C_j|}}$  if  $i \in C_j$   $j = 1 \dots k$   
 properties "h"  
 $H_{ij} = 0$  otherwise  
 $(h_p)$  is the indicator vector -

②  $\mathbb{H}^T \mathbb{H} = \mathbb{I}$   $\mathbb{H}_{(:, p)} = h_p$   
 $h_p^T \mathbb{I} h_p = \frac{\text{cut}(C_p, \bar{C}_p)}{|C_p|} = (\mathbb{H}^T \mathbb{I} \mathbb{L} \mathbb{H})_{pp} \rightarrow$  (product of max matrix of  $\mathbb{H}$  diagonal  $(\mathbb{L}\mathbb{L})$ )  
 $\therefore \text{RatioCut}(C_1, \dots, C_k) = \text{Tr}(\mathbb{H}^T \mathbb{L} \mathbb{H})$   
given :-  $\mathbb{H}^T \mathbb{H} = \mathbb{I}$   
 Trace = RatioCut  
 (similar to what we did for  $k=2$ )

# The RatioCut for arbitrary k

→ for indicator vectors, we impose certain values like  
 $n_{ij} = \frac{1}{\sqrt{|C_j|}}$  if we don't know how to minimize that

# Solving the problem: Relaxation

- Relaxation: Remove the constraint of discrete optimization and accept that the indicator matrix  $\mathbb{H}$  has real values.
- The minimization becomes tractable using Lagrange multipliers. It becomes an eigenvalue-eigenvector problem (see PCA derivation).  
$$\mathbb{H}^T \mathbb{L} \mathbb{H} = \lambda \mathbb{H}$$

solve this eq<sup>z</sup>
- Once it is solved, we need to discretize  $\rightarrow$  k-means using the eigenvectors as coordinates.
- There's no warranty that this solution is similar to the unrelax one!!!  
(non-relaxed)

# The Unnormalized spectral clustering algorithm

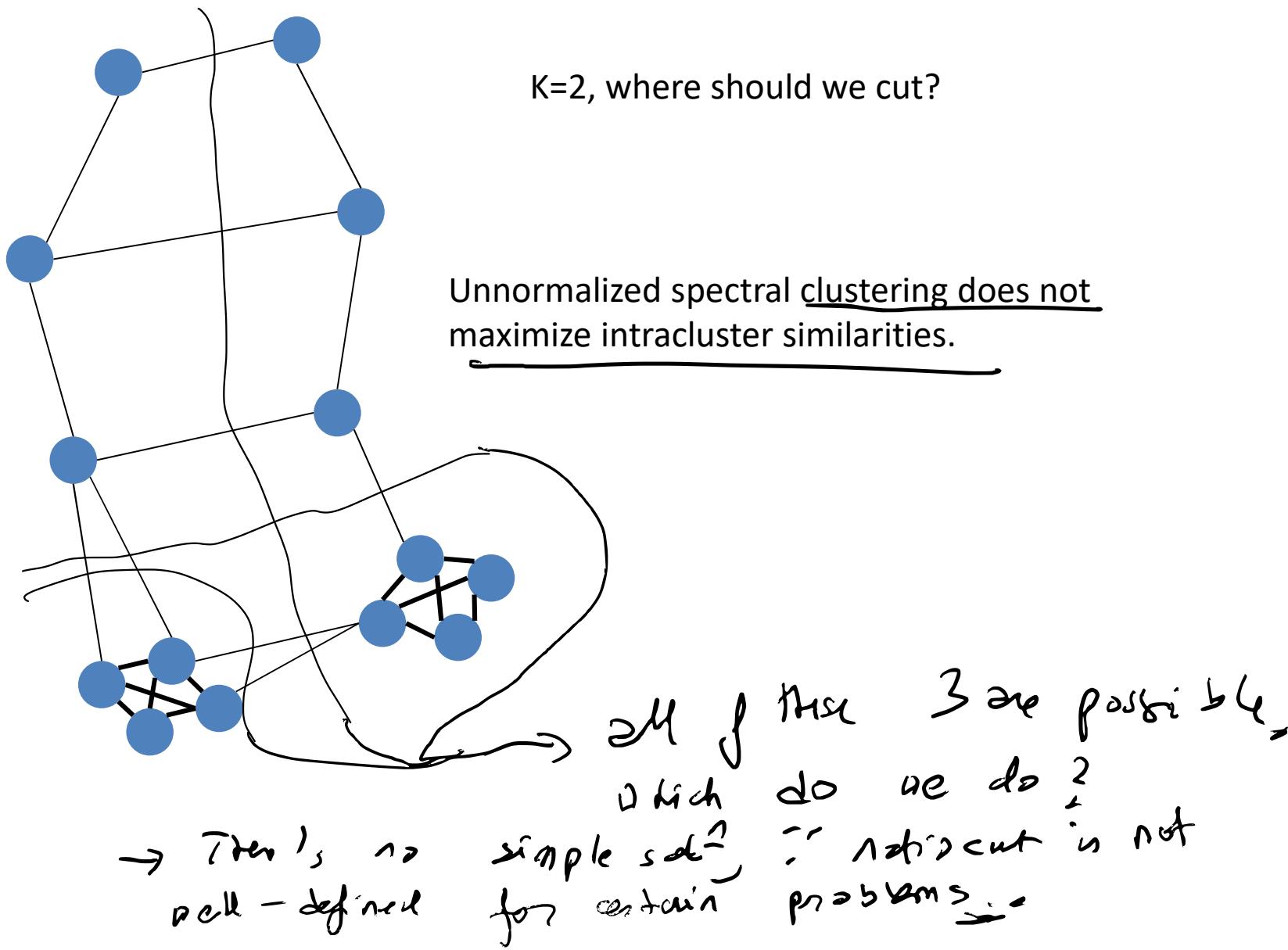
1. Construct the similarity graph \*
2. Compute the Laplacian  $\underline{L} = \mathcal{D} - \mathcal{S}$
3. Compute the first  $k$  eigenvectors  $\underline{L} h = \lambda h$
4. Use the coordinates as input for the  $k$ -means algorithm and obtain the clusters

$$\underline{L} h = \lambda h \Rightarrow \underline{H} \begin{pmatrix} H_{11} & \dots & H_{1k} \\ \vdots & & \vdots \\ H_{n1} & \dots & H_{nk} \end{pmatrix} \xrightarrow{i=1 \dots n} \underline{x}_i = (H_{i1}, \dots, H_{ik})$$

K-means

(new set of centers for k-means)

# A problem related to RatioCut



tries to maximize intracluster similarities

# The Normalized spectral clustering

$$\min N_{\text{cut}} = \frac{\text{Cut}(C, \bar{C})}{\text{Vol}(C)}$$

$$N_{\text{cut}}(C_1 \dots C_k) = \frac{1}{2} \sum \frac{\text{Cut}(C_e, \bar{C}_e)}{\text{Vol}(C_e)}$$

(@)  $\frac{1}{\text{Vol}(C_j)} H_{ij} = 1 \text{ if } i \in C_j ; \text{ otherwise } H_{ij} = 0$

$$\min (\text{Tr}(\mathbb{H}^T \mathbb{L} \mathbb{H})) \text{ given } (\mathbb{H}^T \mathbb{D} \mathbb{H}) = \mathbb{I}$$

minimizing  $N_{\text{cut}}$   
is same as  
minimizing the trace

(B)  $\mathbb{H} = \mathbb{D}^{1/2} \mathbb{L} \mathbb{H}$

$$\min (\text{Tr}(\mathbb{H}^T \mathbb{L} \mathbb{H})) = \min (\text{Tr}(\underbrace{\mathbb{H}^T \mathbb{D}^{-1/2} \mathbb{L} \mathbb{D}^{1/2} \mathbb{H}}_{\mathbb{L}_{\text{sym}}}))$$

cond :-  $(\mathbb{H}^T \mathbb{H}) = \mathbb{I}$

; find problem :-  $\min (\text{Tr}(\underbrace{\mathbb{H}^T \mathbb{L}_{\text{sym}} \mathbb{H}}))$ , given  $(\mathbb{H}^T \mathbb{H}) = \mathbb{I}$

→ Now it is again  
an eigen problem,

& we do same as  
unnormalized algo,  
replace " $\mathbb{L}$ " with " $\mathbb{L}_{\text{sym}}$ "

\* For Ratio Cut the condition was  $\mathbb{H}^T \mathbb{H} = \mathbb{I}$

# Spectral Clustering

(usually it refers to  
normalized  
version)  $\rightarrow$

- Closely related with kernel k-means, under certain conditions, they optimize the same loss.
- With a different transformation, the Laplacian can be transformed in a random walk matrix and related with a Markov state model analysis.
- The value of  $k$  can be inferred from a gap in the spectrum of eigenvalues.

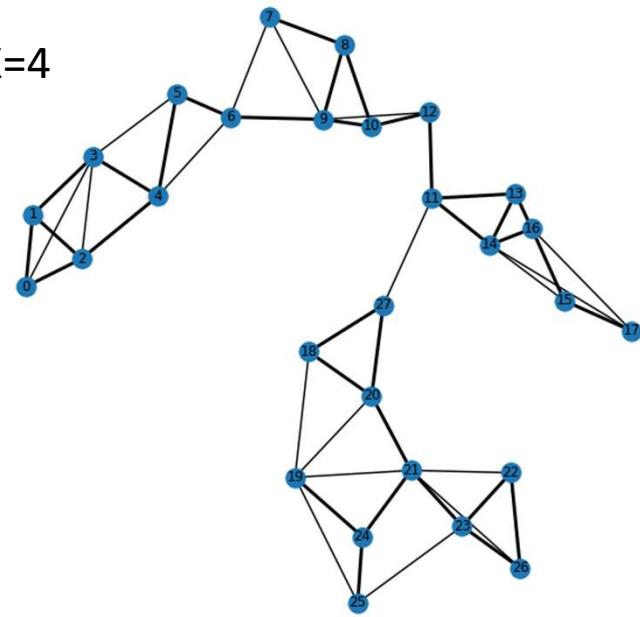
$$\mathcal{L} = \mathcal{D}^{-1} \mathcal{L}$$

# Spectral Clustering problems

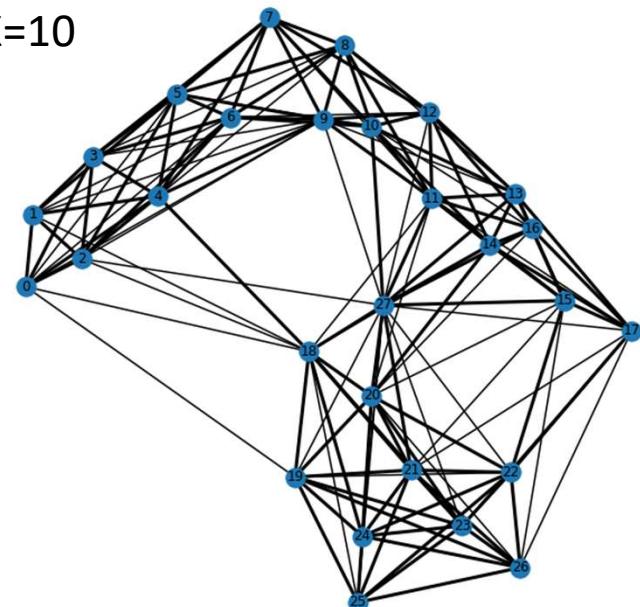
- How to construct the similarity matrix?

K-NN graph

K=4



K=10

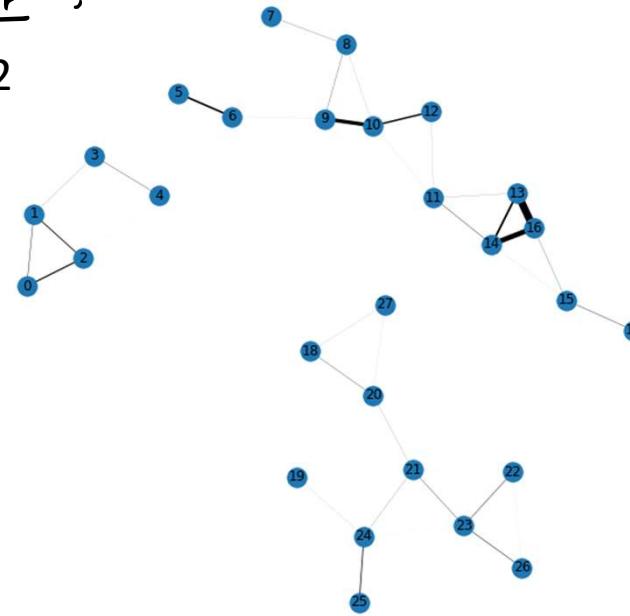


(graphs are completely diff., depending on choice of parameters etc.)

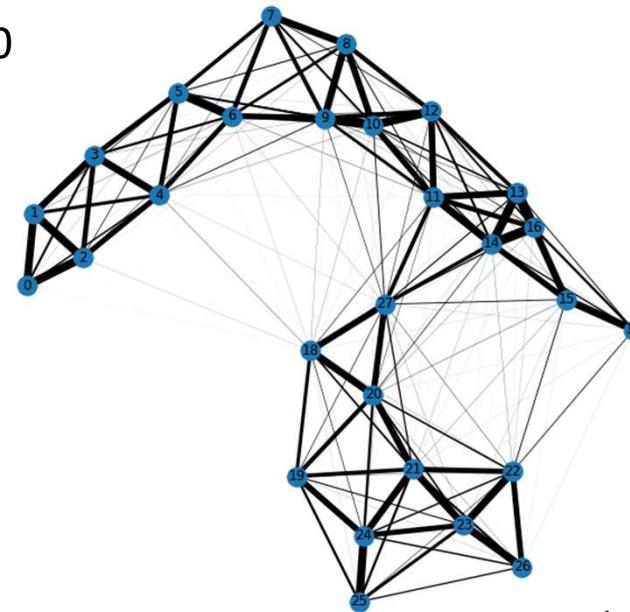
fully connected  $S_{ij} = e^{-\frac{d_{ij}}{\sigma}}$

graph :-

$\sigma=3.2$



$\sigma=1.0$



# Spectral Clustering problems

- How to construct the similarity matrix?
- How to decide the number of clusters in absence of gap?
- Computational Scaling
- Noise sensitivity

$$\xrightarrow[\text{matrix}]{\text{diagonalizing } A} \Theta(n^3)$$

random points added to the graph can cause a mess.

# Expectation-Maximization clustering

Modern clustering algorithms (II)

(EM)) }      Syonyms

## Model based clustering

- Consider your data as a set of realizations of an underlying probability function  $p(x)$ .
- Assume the functional form of  $\underline{p(x)}$  and estimate its parameters by maximum likelihood. (MLE)

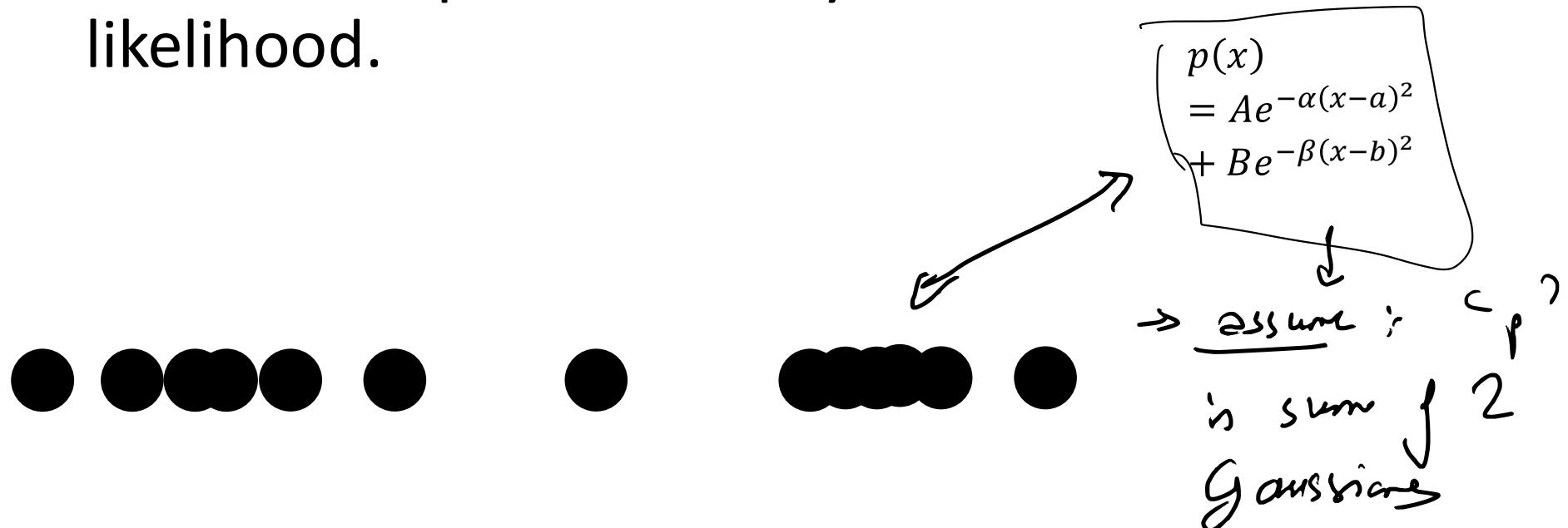
# Model based clustering

- Consider your data as a set of realizations of an underlying probability function  $p(x)$ .
- Assume the functional form of  $p(x)$  and estimate its parameters by maximum likelihood.



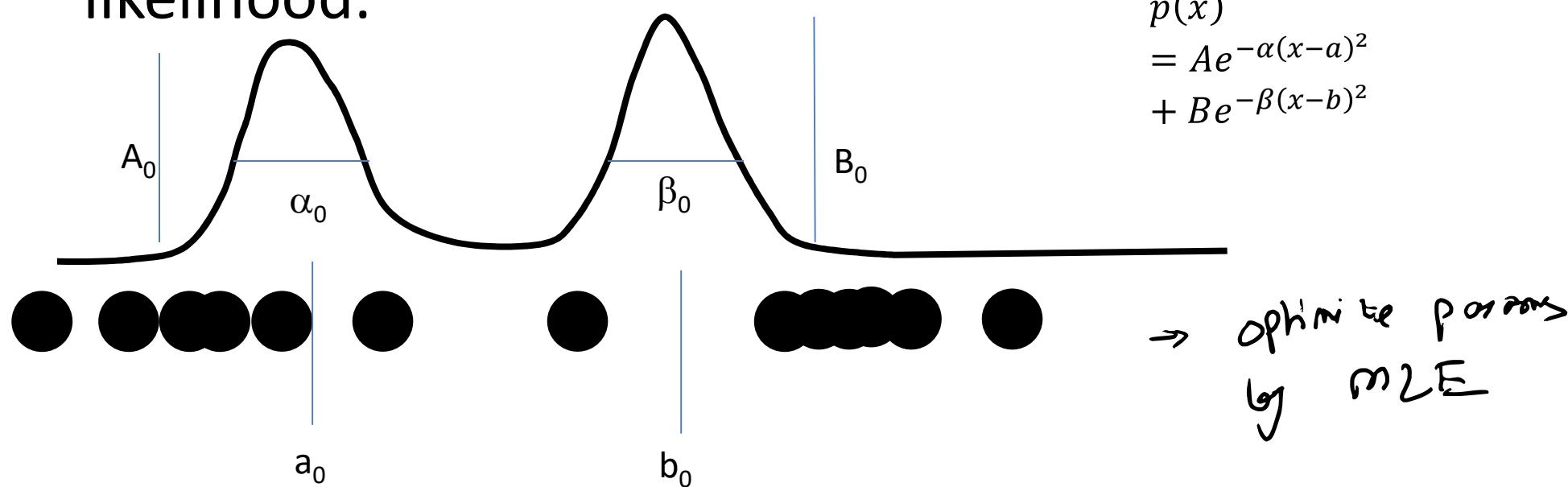
# Model based clustering

- Consider your data as a set of realizations of an underlying probability function  $p(x)$ .
- Assume the functional form of  $p(x)$  and estimate its parameters by maximum likelihood.



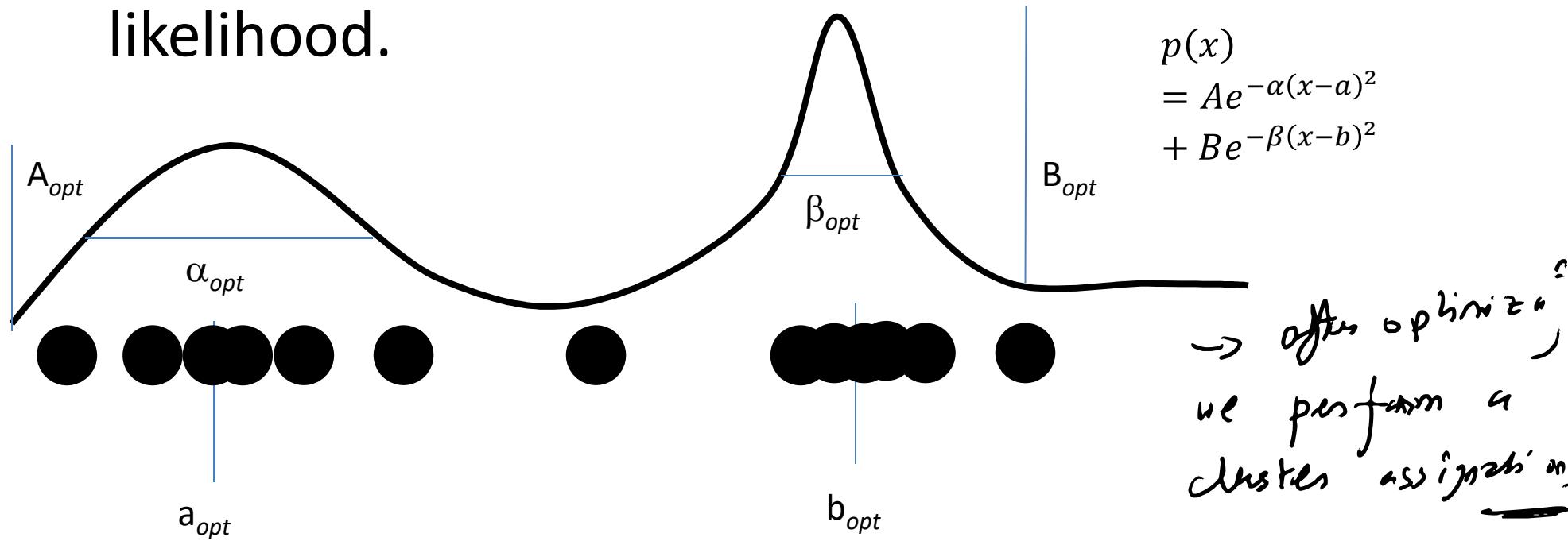
# Model based clustering

- Consider your data as a set of realizations of an underlying probability function  $p(x)$ .
- Assume the functional form of  $p(x)$  and estimate its parameters by maximum likelihood.



# Model based clustering

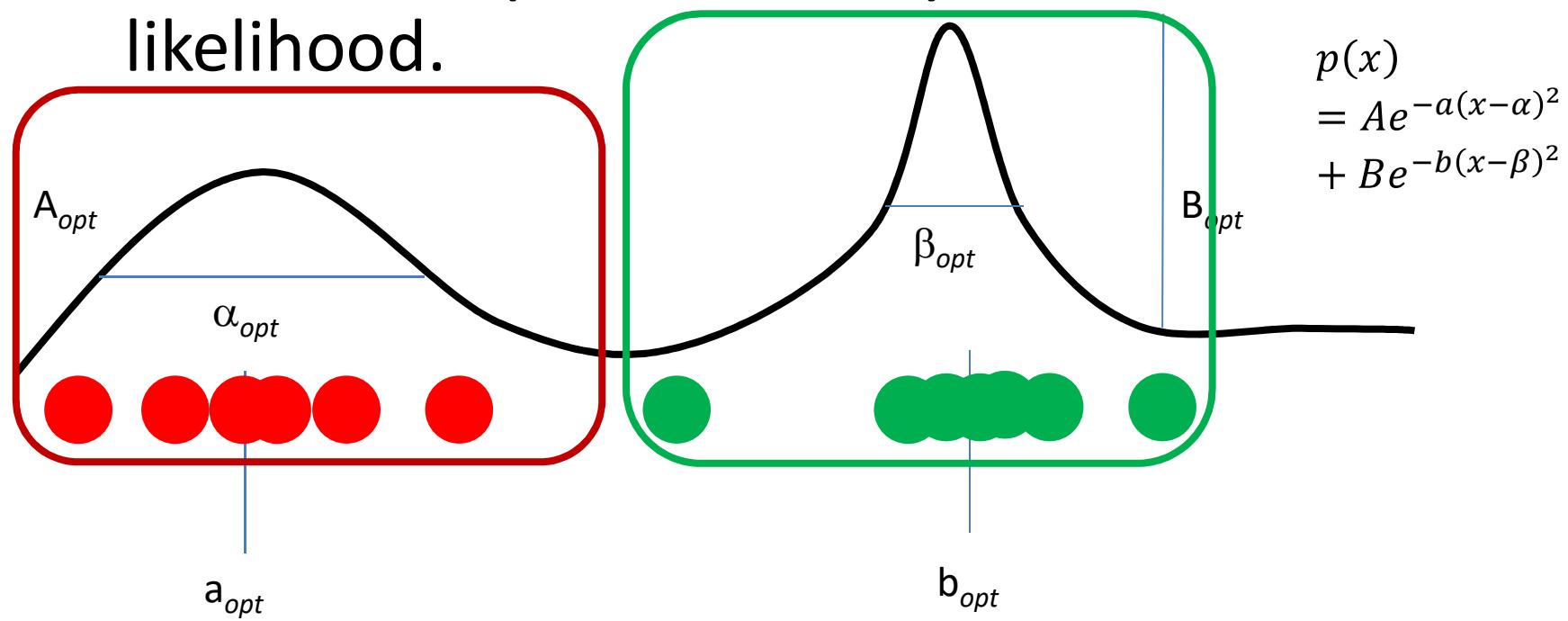
- Consider your data as a set of realizations of an underlying probability function  $p(x)$ .
- Assume the functional form of  $p(x)$  and estimate its parameters by maximum likelihood.



∴ EM algorithm is to estimate the params when they can't be estimated analytically.

## Model based clustering

- Consider your data as a set of realizations of an underlying probability function  $p(x)$ .
- Assume the functional form of  $p(x)$  and estimate its parameters by maximum likelihood.



# Expectation-maximization algorithm

- Iterative procedure to compute the **Maximum Likelihood (ML)** estimate – even in the presence of missing or hidden data
- EM consists of two steps:
  - **Expectation step:** the (missing) data are estimated given the observed data and current estimates of model parameters
  - **Maximization step:** The likelihood function is maximized under the assumption that the (missing) data are known , C.i.e., we recompute the parameters given the prob. distribution  
→ Explained in next slides

# EM algorithm for mixture of Gaussians

- What is a mixture of **K** Gaussians?

$$p(x) = \sum_{k=1}^K \pi_k F(x | \Theta_k)$$

Diagram annotations:

- A curved arrow points from the term  $\pi_k$  to the word "weights".
- A curved arrow points from the term  $F(x | \Theta_k)$  to the words "Gaussian func's".

with

$$\boxed{\sum_{k=1}^K \pi_k = 1}$$

Diagram annotations:

- A curved arrow points from the equation  $\sum_{k=1}^K \pi_k = 1$  to the word "weights".
- The text "weights are normalized too." is written next to the equation.

and **F(x | Θ)** is the Gaussian distribution with parameters **Θ = {μ, Σ}**

# EM algorithm for mixture of Gaussians

- If all points  $x \in X$  are mixtures of  $K$  Gaussians then

$$\mathcal{L}(X) = \prod_{i=1}^n p(x_i) = \prod_{i=1}^n \sum_{k=1}^K \pi_k F(x_i | \Theta_k)$$

- **Goal:** Find  $\pi_1, \dots, \pi_k$  and  $\Theta_1, \dots, \Theta_k$  such that  $\mathcal{L}(X)$  is maximized
- Or,  $\ln(\mathcal{L}(X))$  is maximized: (equivalent)

$$\ln(\mathcal{L}(X))$$

$$L(\Theta) = \sum_{i=1}^n \ln \left\{ \sum_{k=1}^K \pi_k F(x_i | \Theta_k) \right\}$$

$\mathcal{L}(X) = \text{like likelihood}$   
 $L(\Theta) = \text{loss func}^2$

= product  
of  
gaussians  
for  
all  
points  
in  
set  
 $X$   
likelihood  
is  
product  
of  
prob.  
of  
points

# Mixtures of Gaussians -- notes

- Every point  $\mathbf{x}_i$  is probabilistically assigned (generated) to (by) the  $k$ -th Gaussian
- Probability that point  $\mathbf{x}_i$  is generated by the  $k$ -th Gaussian is

$$w_{ik} = \frac{\pi_k F(x_i | \Theta_k)}{\sum_{j=1}^K \pi_j F(x_i | \Theta_j)}$$

# Mixtures of Gaussians -- notes

- Every Gaussian (cluster)  $C_k$  has an effective number of points assigned to it  $N_k$

$$N_k = \sum_{i=1}^n w_{ik}$$

*effective # of pts comes from sum of weights*

- With mean

$$\mu_k = \frac{1}{N_k} \sum_{i=1}^n w_{ik} \underline{x_i}$$

compute  
mean

- And variance

$$\Sigma_k = \frac{1}{N_k} \sum_{i=1}^n w_{ik} (x_i - \mu_k) x_i (x_i - \mu_k)^T$$

covariance  
matrix

# EM for Gaussian Mixtures

- Initialize the means  $\mu_k$ , variances  $\Sigma_k$  ( $\Theta_k = (\mu_k, \Sigma_k)$ ) and mixing coefficients  $\pi_k$ , and evaluate the initial value of the loglikelihood

(C can be n random or from k means)

- **Expectation step:** Evaluate weights

$$w_{ik} = \frac{\pi_k F(x_i | \Theta_k)}{\sum_{j=1}^K \pi_j F(x_i | \Theta_j)}$$

how likely is it that  
each point belongs  
to a given cluster.

# EM for Gaussian Mixtures

- **Maximization step:** Re-evaluate parameters

(Compute new params for the Gaussians)

$$\mu_k^{new} = \frac{1}{N_k} \sum_{i=1}^n w_{ik} x_i$$

$$\Sigma_k^{new} = \frac{1}{N_k} \sum_{i=1}^n w_{ik} (x_i - \mu_k^{new}) x_i (x_i - \mu_k^{new})^T$$

$$\pi_k^{new} = \frac{N_k}{N}$$

- Evaluate  $L(\Theta^{new})$  and stop if converged

# EM characteristics

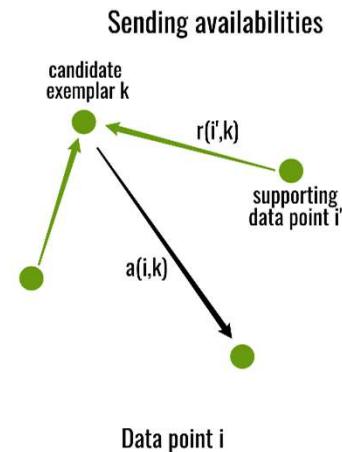
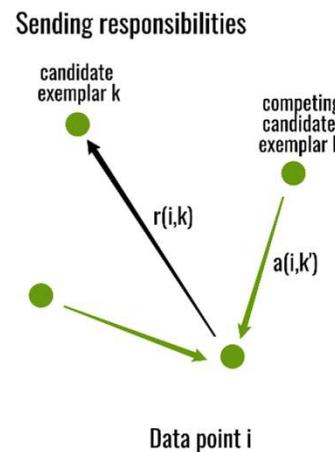
- Notice the similarity between EM for Normal mixtures and K-means; The expectation step is the assignment, while the maximization step is the update of the centers.
- Contrary to what happens in the case of fuzzy clustering, the weights have a real probabilistic interpretation.
  - If the model ( $k$ , functional form) is not realistic, neither will the results.
  - It is not guaranteed to reach the global optimum.
  - $K$  can be inferred with (bayesian) model selection (Dirichlet process).
    - "prior" is imposed to often " $\pi$ ".
    - So we need to report many times like in K-means.

# Affinity propagation

Modern clustering algorithms (III)

# Affinity Propagation. Concepts

- AffinityPropagation creates clusters by sending messages between pairs of samples until convergence.
- Two categories of messages:
  - **Responsibility**
  - **Availability**



# Affinity Propagation. Responsibility

Accumulated evidence that sample  $k$  should be the exemplar for sample  $i$ .

$$r(i, k) \leftarrow \underline{S_{ik}} - \max_{k' \neq k} [\underline{S_{ik'}} + \underline{a(i, k')}]$$

↓ at = points that are representatives of the clusters

→ compare evidence that a point is exemplar of another.

(This is no formula)

# Affinity Propagation. Availability

Accumulated evidence that sample  $i$  should choose sample  $k$  to be its exemplar

$$a(i, k) \leftarrow \min \left[ 0, r(k, k) + \underbrace{\sum_{i' \neq \{k, i\}} \max\{0, r(i', k)\}}_{\text{for itself}} \right]$$
$$a(k, k) \leftarrow \underbrace{\sum_{i' \neq \{k\}} \max\{0, r(i', k)\}}_{\text{for itself}}$$

Note that it considers the values for all other samples

# Affinity Propagation. The algorithm

- $S_{ik}$  is the similarity, usually is taken  $S_{ik} = -\|x_i - x_k\|^2$  *(-ve ∫ dist b/w<sup>2</sup> dat, p.k. ⇒ common is affinity prop.)*
- $S_{kk}$ , known as preference, plays an important role. A value close to the minimum possible similarity produces fewer classes, while a value close to or larger than the maximum possible similarity produces many classes. It is typically initialized to the median similarity of all pairs of inputs.

# Affinity Propagation. The algorithm

- The algorithm starts by setting all the values of  $r$  and  $a$  to zero. And iterates by first updating responsibilities and then availabilities until a number of iterations is reached.
- So, after the first iteration:

$$r(i, k) \leftarrow S_{ik} - \max_{k' \neq k} [S_{ik'}] \quad \text{available}$$

$$a(i, k) \leftarrow \min[0, r(k, k) + \sum_{i' \neq \{k, i\}} \max\{0, r(i', k)\}]$$

- $r(i, k)_t \leftarrow S_{ik} - \max_{k' \neq k} [S_{ik'} + a(i, k')]$  decreases when the availability for other points increases!

if  $(i)$  &  $(k)$  are most similar,  $r > 0$

if  $(i)$  &  $(k')$  are most similar,  $r < 0$

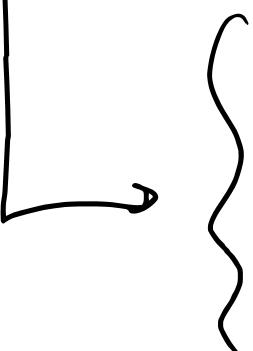
# Affinity Propagation. The algorithm

- The exemplars are extracted from the final matrices as those whose 'responsibility + availability' for themselves is positive ( $r(k, k) + a(k, k) > 0$ ).
- The element  $i$  is assigned to the same cluster as the exemplar  $l$  for which  $r(i, l)$  is maximum.

# Affinity Propagation. The algorithm

- The algorithm described in this way is usually not stable numerically.
- A damping factor is introduced to smooth the transitions:

$$\begin{cases} r(i, k)_{t+1} = \lambda \underbrace{r(i, k)_t}_{\text{one or re-compute } f(a)} + (1 - \lambda) \overline{r(i, k)_{t+1}} \\ a(i, k)_{t+1} = \lambda \underbrace{a(i, k)_t}_{\text{"damp" the change by a factor}} + (1 - \lambda) \overline{a(i, k)_{t+1}} \end{cases}$$

 → one or re-compute  $f(a)$ , we "damp" the change by a factor  $\rightarrow$  that in transitions are smoother,

# Affinity Propagation. Prices and caveats

- Two parameters: Preference and dumping factor.
- The preference parameter strongly affects the results and setting it is not so clear.
- The method is quite weight computationally speaking, with a complexity that is  $\mathcal{O}(\underline{T}\underline{N}^2)$ .
- It trends to generate convex clusters.

# Next lecture...

Metrics for validation

- K-means
  - Fuzzy c-means
  - Agglomerative hierarchical
  - Spectral Clustering
  - Gaussian Mixture  $(\text{Expectation} - \text{Maximization})$
  - Affinity Propagation
  - Density Based Algorithms:
    - DBSCAN
    - Density Peaks
    - Mean-Shift clustering
- 