

- A different family of approaches builds a graph on data and attempt to learn clusters from the structure in the graph.
- Let's think of a graph G as a collection of vertices V and edges between the vertices E : $G = (V, E)$.
- We will associate V with the data points and E will "weights" between points.
- E will often be denoted as a $|V| \times |V|$ matrix W , where $W_{ij} \in [0, 1]$, $i, j = 1, \dots, |V|$.
 $\rightarrow \{x_i\}_{i=1}^n \subset \mathbb{R}^D$
 $W_{ij} \sim 0 \Rightarrow$ weak connection between x_i, x_j
 $W_{ij} \sim 1 \Rightarrow$ strong connection between x_i, x_j .

ex: Fix some integer $k_{NN} \geq 1$. Let $W_{ij} = \begin{cases} 1, & x_i \text{ is among the } k_{NN}\text{-nearest neighbors of } x_j \\ 0, & \text{else} \end{cases}$

This is a k_{NN} -nearest neighbors graph, where x_i a k_{NN} -nearest neighbor of x_j means that when all pairwise distances $\{\|x_\ell - x_j\|_2\}_{\ell=1}^n$ are computed, $\|x_i - x_j\|_2$ is among the k_{NN} smallest.

Remark: If $k_{NN} \ll n$, this is a sparse graph. Computationally

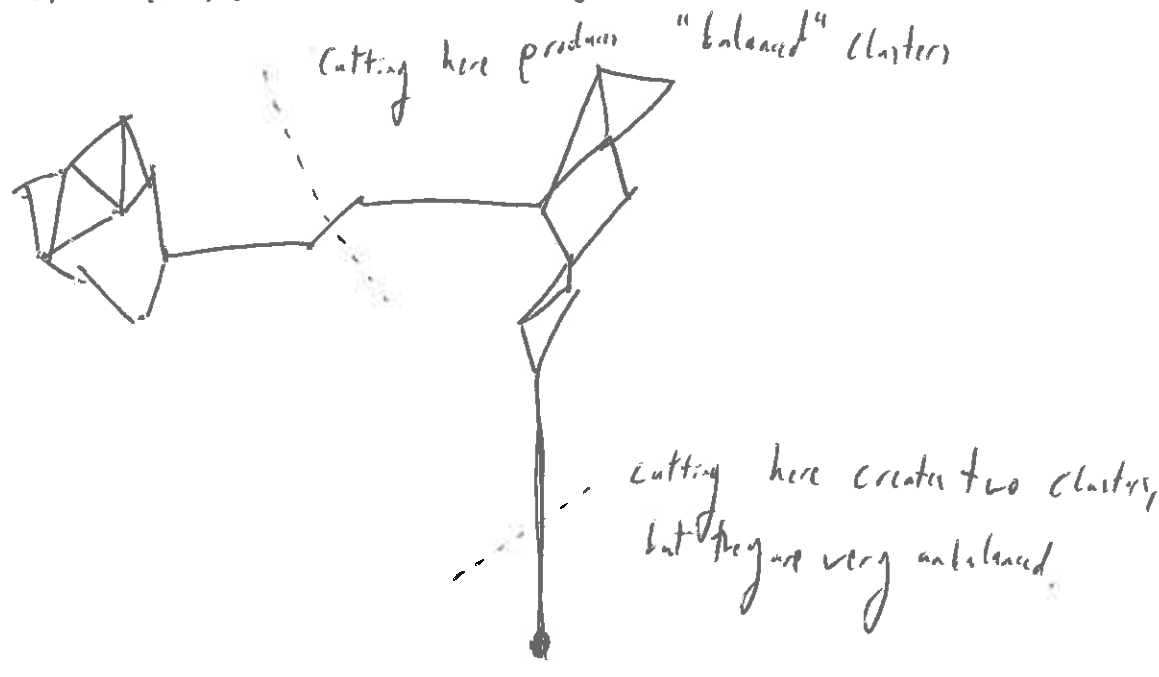
advantageous in some cases.

ex: $w_{ij} = e^{-\frac{\|x_i - x_j\|_2^2}{2\sigma^2}}$, $\sigma > 0$ fixed scale parameter. This connects

all points since $e^{-\frac{\|x_i - x_j\|_2^2}{2\sigma^2}} > 0$ no matter what $\frac{\|x_i - x_j\|_2^2}{2\sigma^2}$ is, but

it may be very small.

Given these weights on the graph, we can try to break the graph into ^{two} pieces that "cut" as little structure as possible:



Attempt 1: Partition so that the weight of edges broken is as small as possible.

$Cut(C_1, C_2) = \sum_{\substack{x_i \in C_1 \\ x_j \in C_2}} w_{ij}$. Choose clusters to minimize this cut:

$(C_1^*, C_2^*) = \underset{(C_1, C_2)}{\operatorname{argmin}} Cut(C_1, C_2) = \underset{(C_1, C_2)}{\operatorname{argmin}} \sum_{\substack{x_i \in C_1 \\ x_j \in C_2}} w_{ij}$

• Seems reasonable... but could end up breaking a small number of edges connecting outliers to the main data bulk. (3)

Attempt 2: Partition so that the weight of the edges broken is small and also so that the "size" of the clusters is large.

$$N_{cut}(C_1, C_2) = \frac{\sum_{x_i \in C_1} \sum_{x_j \in C_2} w_{ij}}{\sum_{x_i \in C_1} \sum_{x_j \in C_1} w_{ij}} + \frac{\sum_{x_i \in C_1} \sum_{x_j \in C_2} w_{ij}}{\sum_{x_i \in C_2} \sum_{x_j \in C_2} w_{ij}}$$

$$= \sum_{k=1}^2 \frac{\sum_{x_i \in C_k} \sum_{x_j \in C_k} w_{ij}}{\sum_{x_i, x_j \in C_k} w_{ij}}$$

Setting $vol(C_k) = \sum_{x_i \in C_k} \sum_{x_j \in C_k} w_{ij}$, this simplifies to

$$N_{cut}(C_1, C_2) = \sum_{k=1}^2 \left(\frac{\sum_{x_i \in C_k} \sum_{x_j \in C_k} w_{ij}}{vol(C_k)} \right)$$

- Compared to Attempt 1, this is normalized by the volume factor $vol(C_k)$.
- In particular, if $vol(C_k)$ is very small, N_{cut} blows up... penalizing small clusters, thus enforcing balance.

$$(C_1^*, C_2^*) = \underset{(C_1, C_2)}{\operatorname{arg\,min}} \operatorname{Ncut}(C_1, C_2)$$

$$= \underset{(C_1, C_2)}{\operatorname{arg\,min}} \sum_{k=1}^2 \frac{\sum_{\substack{x_i \in C_k \\ x_j \notin C_k}} w_{ij}}{\operatorname{vol}(C_k)} \quad \textcircled{*}$$

Questions: 1.) How to generalize to more than 2 clusters.
 2.) How to compute.

Answers: 1.) Easy: if you want K clusters, extend $\textcircled{*}$ to sum over K partition elements:

$$(C_1^*, C_2^*, \dots, C_K^*) = \underset{(C_1, C_2, \dots, C_K)}{\operatorname{arg\,min}} \sum_{k=1}^K \frac{(\sum_{\substack{x_i \in C_k \\ x_j \notin C_k}} w_{ij})}{\operatorname{vol}(C_k)}.$$

2.) Hard: Computing the exact solution to $\textcircled{*}$ is NP-hard. It's a combinatorial nightmare, and one can prove there is no efficient algorithm. Basically, you have to check all possibilities in the worst case, which means exponentially many in n ~~same~~ computation. Bad!

• Instead of calculating $\textcircled{*}$ exactly, we can use linear algebra to relax $\textcircled{*}$.

• Indeed, let us consider a partition (C_1, C_2) . Let f_{ij}^C be defined as:

$$\bar{C} = V \mathbf{1} C \quad \text{''} \quad (C, \bar{C})$$

$$f^c_i = \begin{cases} \sqrt{\frac{\text{vol}(\bar{c})}{\text{vol}(c)}} & , x_i \in c \\ -\sqrt{\frac{\text{vol}(c)}{\text{vol}(\bar{c})}} & , x_i \in \bar{c} \end{cases}$$

5

This is + for points in c , - for points in \bar{c} , and moreover is closely related to the degree matrix:

$$D_{ij} = \begin{cases} \sum_{x_j \in V} w_{ij} & , i = l \\ 0 & , \text{else} \end{cases}$$

$$D = \begin{pmatrix} \sum_j w_{1j} & & & 0 \\ & \sum_j w_{2j} & & \\ & & \ddots & \\ 0 & & & \sum_j w_{nj} \end{pmatrix} \quad \text{Then } (Df^c)^T \mathbf{1} = 0,$$

\uparrow column vector

where $\mathbf{1} = (1, \dots, 1)^T$.

Moreover, $(f^c)^T D f^c = \text{vol}(V) = \sum_{x_i, x_j} w_{ij}$.

Most importantly, f^c is intimately related to $N_{\text{cut}}(c, \bar{c})$ through the graph Laplacian $L = D - W$. Then $f^c L f^c = \text{vol}(V) N_{\text{cut}}(c, \bar{c})$.

-So, instead of handling \mathbb{A} directly, we can use this Laplacian formulation:

$$c^* = \underset{c}{\text{argmin}} (f^c)^T L f^c \quad \text{s.t.} \quad \begin{aligned} &1.) (Df^c)^T \mathbb{1} = 0 \\ &2.) (f^c)^T D f^c = \text{vol}(V). \end{aligned}$$

• Great... so what? Still hard to solve. The problem is f^c has special form.

$$\text{Relat} = \underset{f \in \mathbb{R}^n}{\text{argmin}} f^T L f \quad \text{s.t.} \quad \begin{aligned} &1.) (Df)^T \mathbb{1} = 0 \\ &2.) f^T D f = \text{vol}(V). \end{aligned}$$

↓
lose some accuracy, gain ease of computation.

This can be attacked with linear algebra... next time.