# SPECTRAL CLUSTERING

<u>Notation</u> : Let $G = (V, E)$ - undirected, weighted graph

$$V = \{v_1, \ldots, v_n\}, \quad v_i \overset{w_{ij}}{\bullet\!-\!-\!\bullet} v_j \quad w_{ij} \geq 0$$

- adjacency matrix : $\quad W = (w_{ij})_{i,j=1,\ldots,n} \quad [W = W']$

- degree matrix : $\quad d_i = \sum_{j=1}^{n} w_{ij}, \quad D = diag(d_1, \ldots, d_n)$

<u>Similarity Graphs</u> : Transform data $\{x_1, \ldots, x_n\}$ with pairwise
   distances $d_{ij}$ into graph

① $\epsilon$ - neighbourhood graph

   -- connect all points whose pairwise ~~edges~~
   distances are less than $\epsilon$.

   - since $w_{ij} \leq \epsilon$, can be considered un-weighted

② $k$ - nearest - neighbour graph :

   - connect $v_i$ to $v_j$ if $v_i$ is among $k$-nearest
   neighbours of $v_j$ and $v_j$ is among $k$-nearest
   neighbors of $v_i$.

   - weight the edges by similarity of $v_i, v_j$

③ fully connected graph :

   - connect all points with positive "similarity"

   - eg : $S_{ij} = S(x_i, x_j) = \exp\left(\dfrac{-\|x_i - x_j\|^2}{2\sigma^2}\right)$

# GRAPH LAPLACIANS

① Unnormalized graph laplacian :

$$L = D - W$$

**Proposition 1 :** ⓐ $f'Lf = \frac{1}{2} \sum_{i,j} w_{ij} (f_i - f_j)^2 \quad f \in \mathbb{R}^n$

ⓑ $L$ – symmetric, p.s.d

ⓒ $\lambda_{min} = 0$, $u_{min} = \mathbb{1} \in \mathbb{R}^n$

ⓓ $0 \le \lambda_{min} = \lambda_1 \le \lambda_2 \le \dots \le \lambda_n$

**Proof :** ⓐ $x'Lx = x'Dx - x'Wx$

$$= \sum_i d_i x_i^2 - \sum_{i,j} x_i x_j w_{ij}$$

$$= \frac{1}{2} \left( \sum_i d_i x_i^2 - 2 \sum_{i,j} x_i x_j w_{ij} + \sum_j d_j x_j^2 \right)$$

$$= \frac{1}{2} \left( \sum_i \sum_j w_{ij} x_i^2 - 2 \sum_{i,j} x_i x_j w_{ij} + \sum_j \sum_i w_{ji} x_j^2 \right)$$

$$= \frac{1}{2} \sum_{i,j} w_{ij} (x_i - x_j)^2$$

ⓑ – $D, W$ are symmetric ⟹ $L$ is symmetric

– $x'Lx \ge 0 \quad \forall x$ ⟹ $L$ is p.s.d

ⓒ – choosing $x = \mathbb{1}$ ⟹ $\lambda_{min} = 0$

ⓓ follows from ⓐ, ⓑ, ⓒ : ▢

Proposition 2: The multiplicity $k$ of eigenvalue $0$ of $L$ equals the number of connected components $A_1, ..., A_k$. The eigenspace of $0$-eval is spanned by vectors $\mathbb{1}_{A_1}, ... \mathbb{1}_{A_k}$

Proof : I: let $k = 1$; $\Rightarrow$ graph is connected. Let $x$ be eigenvector with eigenvalue $0$. then

$$x'Lx = \frac{1}{2} \sum w_{ij} (x_i - x_j)^2 = 0$$

$\Rightarrow$ $x$ needs to be eqval on all nodes which can be connected by a path in $G$. $\Rightarrow$ $x_i = \mathbb{1}$ $\forall i$

II w.l.o.g, let $L = \begin{pmatrix} L_1 & & 0 \\ & L_2 & \\ 0 & & \ddots & \\ & & & L_k \end{pmatrix}$

- $L_i$ are graph laplacians of $A_i$
- eigenvectors of $L$ = eigenvectors of $L_i$ with $0$'s added appropriately
- eigenvalues of $L_i$ = e-values of $L$. $\Rightarrow$ there are $k$ eigenvalues $0$, and the eigenvector in $\mathbb{1}_{A_i}$ [ $\mathbb{1}_A = (x_1, ..., x_n)^T$

$\Leftrightarrow$ $x_i = 1$ if $v_i \in A$ ]

② Normalized Graph Laplacians:

$$L_{sym} = D^{-1/2} L D^{-1/2} = I - D^{-1/2} W D^{-1/2}$$

$$L_{rw} = D^{-1} L = I - D^{-1} W$$

③

**Proposition 3 :** The normalized laplacians satisfy following properties.

(a) $x' L_{sym} x = \frac{1}{2} \sum_{i,j} w_{ij} \left( \frac{x_i}{\sqrt{d_i}} - \frac{x_j}{\sqrt{d_j}} \right)^2$

(b) $\lambda$ is an eigenvalue of $L_{rw}$ with e-vector $u$ if and only if $\lambda$ is an eigenvalue of $L_{sym}$ with e-vector $w = D^{-1/2} u$

(c) $\lambda$ is an eigenvalue of $L_{rw}$ with e-vector $u$ iff $\lambda$ and $u$ solve the generalised eigenproblem $L u = \lambda D u$

(d) $0$ is an eigenvalue of $L_{rw}$ with the constant $\mathbb{1}$ as eigenvector. $0$ is an eigenvalue of $L_{sym}$ with e-vector $D^{-1/2} \mathbb{1}$

⊗

**Proposition 4 :** The multiplicity $k$ of eigenvalue $0$ of both $L_{rw}$ and $L_{sym}$ equals the number of connected components $A_1, \ldots, A_k$ in $G$. For $L_{rw}$, eigenspace of $0$ is spanned by $\mathbb{1}_{A_i}$ and for $L_{sym}$, eigenspace of $0$ is spanned by $D^{1/2} \mathbb{1}_{A_i}$

Proof of Prop 3, 4 are similar to that of Props. 1, 2 ⊗

# SPECTRAL CLUSTERING ALGORITHM

## Un-normalized spectral clustering

Input : Similarity matrix $S \in \mathbb{R}^{n \times n}$ # clusters $K$

1. Construct similarity graph with adj. matrix $W$
2. Compute un-normalized Laplacian $L$
3. Compute top-$k$-eigenvectors of $L$ as $V = [v_1, \ldots, v_k] \in \mathbb{R}^{n \times k}$
4. Let $y_i$ $(i=1, \ldots, n)$ be $i^{th}$ row of $U$
5. cluster $y_i$ into clusters $C_1, \ldots, C_k$ using $k$-means.

Output : Clusters $A_1, \ldots, A_k$ with $A_i = \{j \mid y_j \in C_i\}$

## Justification OF S.C ALGORITHM [Using Matrix Perturbation]

- ideally, "inter-cluster similarity" is $0$
  $\Rightarrow$ $y_i \in \mathbb{R}^k$ are of the form $[0, 0, \ldots, \overset{j}{1}, 0 \ldots, 0]$
  where $j$ is s.t $y_i \in A_j$ $\Rightarrow$ $x_i \in A_j$

- in a nearly ideal case, "inter-cluster similarity" is $\epsilon$
  $y_i$ will be of form $[0, \ldots 0, 1, 0, \ldots, 0] + \epsilon$

- using davis kahan if $\check{A} = A + H$, $S_1 \subseteq \mathbb{R}$ interval
  $\sigma_{S_1}(A)$ : set of all eig-vals contained in $S_1$ and $V_1$
  be eig. space corresponding to $\sigma_{S_1}(A)$ and same for
  $\check{A}, \check{V_1}$. then
  $$\| \sin \theta(v_1, \check{v}_1) \| \leq \frac{\| H \|}{\delta}$$
  $\delta = \min\{ |\lambda - s| ; \lambda - \text{eigval of } A, \lambda \notin S_1, s \in S_1\}$