

Objective: Given a network, made of nodes and observed edges, the aim is to find
unknown groupings of nodes 17/5/21

↓
clusters, community, block
↓
community detection method

We are going to look at two methods:

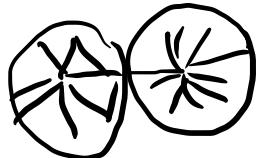
- ① Community detection via hierarchical clustering
- ② Spectral analysis → K-means

COMMUNITY DETECTION VIA HIERARCHICAL CLUSTERING

We need to define a measure of quality of a cluster \rightsquigarrow Think of $W(F_i)$ in K-means

Let $(C_1, \dots, C_K) = C$ a partition of the nodes in a network.

Let $f_{ij}(c)$ denote the fraction of edges in the original network that connect nodes in cluster C_i with nodes in cluster C_j .



In particular, $f_{kk}(c)$ is the fraction of edges within cluster C_k .

MODULARITY
(K fixed)

$$M(c) = \sum_{j=1}^K \left(\underbrace{f_{jj}(c)}_{\text{percentage of edges in } C_j} - \underbrace{f_{jj}^*(c)}_{\text{percentage of edges under random assignment (standardization)}} \right)^2$$

Large values of modularity suggest non-trivial grouping structures

Hierarchical clustering can be based on modularity:

- ① AGGLOMERATIVE: Successive coarsening of partitions through the process of merging.
At each step, consider the merge that has the highest modularity.
- ② DIVISIVE: Successive refinement of partitions through the process of splitting.
At each step, consider the split with the highest modularity.

Implementations would then return the optimum K , ie the partition with the highest modularity overall.

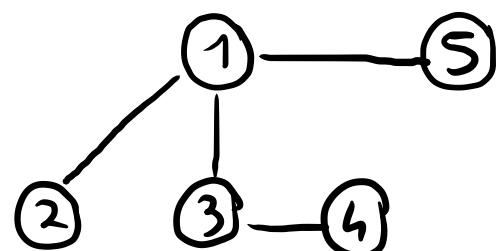
COMMUNITY DETECTION VIA SPECTRAL CLUSTERING

The connectivity of a graph (network) is related to the eigenanalysis of its associated Laplacian matrix

Def Given a network, one can define its adjacency matrix:

$$A = (a_{ij}) \text{ where } a_{ij} = \begin{cases} 1 & \text{if node } i \text{ and } j \text{ are connected} \\ 0 & \text{otherwise} \end{cases}$$

Ex:



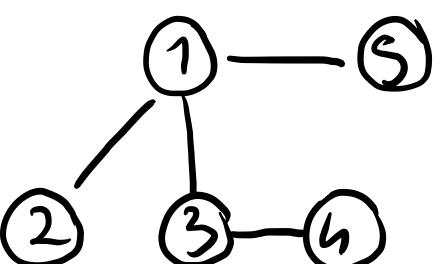
$$A = \begin{pmatrix} 0 & 1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \end{pmatrix}$$

Def Laplacian matrix is defined as

$$L = D - A$$

where A is the adjacency matrix and D is a diagonal matrix with node degrees in the diagonal.

Ex



$$A = \begin{bmatrix} 0 & 1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$D = \begin{bmatrix} 3 & & & & \\ & 1 & & & \\ & & 2 & & \\ & & & 1 & \\ & & & & 1 \end{bmatrix}$$

So the Laplacian is

$$L = \begin{bmatrix} 3 & -1 & -1 & 0 & -1 \\ -1 & 1 & 0 & 0 & 0 \\ -1 & 0 & 2 & -1 & 0 \\ 0 & 0 & -1 & 1 & 0 \\ -1 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Row and column
sums are zero
since $d_i = \sum_{j=1}^r a_{ij}$

L has some nice mathematical properties:
 $p \times p$ number of nodes

- ① For every vector $\underline{v} \in \mathbb{R}^p$,

$$\underline{v}^t L \underline{v} = \frac{1}{2} \sum_{i,j=1}^p a_{ij} (v_i - v_j)^2$$

Proof:

$$\begin{aligned} \underline{v}^t L \underline{v} &= \underline{v}^t (D - A) \underline{v} = \underline{v}^t D \underline{v} - \underline{v}^t A \underline{v} = \sum_{i=1}^p d_i v_i^2 - \sum_{i,j=1}^p a_{ij} v_i v_j \\ &\quad \downarrow \\ &\quad (D)_{ii} \end{aligned}$$

$$= \frac{1}{2} \left[\sum_{i=1}^p d_i v_i^2 + \sum_{j=1}^p d_j v_j^2 - 2 \sum_{i,j=1}^p a_{ij} v_i v_j \right]$$

$$= \frac{1}{2} \left[\sum_{i,j=1}^p a_{ij} v_i^2 - 2 \sum_{i,j=1}^p a_{ij} v_i v_j + \sum_{j=1}^p a_{jj} v_j^2 \right]$$

$$d_i = \sum_{j=1}^p a_{ij}$$

$$= \frac{1}{2} \sum_{i,j=1}^p a_{ij} (v_i - v_j)^2$$

•

② L is symmetric and positive semi-definite

$$l_{ij} = d_{ij} - a_{ij} \neq d_{ji} - a_{ji} = l_{ji} \quad \text{symmetric}$$

D, A
symmetric

$$\text{For any } v \in \mathbb{R}^P, \quad v^T L v = \sum_{i,j} a_{ij} (v_i - v_j)^2 \geq 0 \quad \text{positive semi-definite}$$

④

③ L has P non-negative real-valued eigenvalues

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

④ The smallest eigenvalue is zero and the corresponding eigenvector is the constant one vector.

Take $v = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$, then $Lv = \lambda v$ with $\lambda = 0$

Indeed

$$L \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} = 0 \cdot \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

row sums
are zero in L

This property can be extended to the case of weighted graphs

$$d_i = \sum_{j=1}^P w_{ij} \text{ with } w_{ij} \geq 0$$

$$L = D - W$$

LAPLACIAN AND CONNECTIVITY

Def : A connected component (or just a component) of a graph is a subgraph in which any two vertices are connected by a path (ie a sequence of connected nodes)

Prop The number of zero eigenvalues of the Laplacian matrix is equal to the number of connected components in the graph.

Proof Let us consider the case of 2 connected components.

Then I can rearrange the adjacency matrix like this:

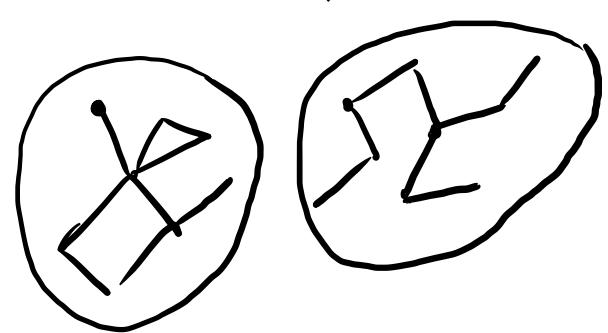
$$A = \begin{pmatrix} A_1 & 0 \\ 0 & A_2 \end{pmatrix} \left. \begin{array}{l} \text{\{} \# nodes in component 1 \\ \{} \# nodes in component 2 \end{array} \right\}$$

Then consider

$$\underline{v}_1 = \begin{pmatrix} 1 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \left. \begin{array}{l} \text{\{} \# nodes in component 1 \\ \{} \# nodes in component 2 \end{array} \right\}$$

$$\underline{v}_2 = \begin{pmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

These two are eigenvectors with eigenvalue zero.



Two zero eigenvalues

with K components:

$$A := \begin{pmatrix} A_1 & & & \\ & A_2 & & \\ & & \ddots & \\ & & & A_K \end{pmatrix}$$

$$v_1 = \begin{pmatrix} 1 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \left. \right\} \text{\# nodes in component 1}$$

$$v_2 = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \{ \text{# nodes in component 2} \}$$

$$v_K = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \hline 1 \\ \vdots \\ 1 \end{pmatrix} \xrightarrow{\text{nodes in component } K}$$

v_1, \dots, v_k are eigenvectors with zero eigenvalues.

EXAMPLE (K=3)

Consider the matrix of eigenvectors:
with zero eigenvalues

\underline{v}_1	\underline{v}_2	\underline{v}_3	
1	0	0	
1	0	0	
1	0	1	
0	1	1	
...	
0	0	1	

equal rows

In practice, networks will not have completely disconnected components.

SPECTRAL CLUSTERING works by applying K-means to the $p \times K$ matrix of eigenvectors corresponding to the smallest eigenvalues.