

Today: community detection (topic of HW1, Wed.)

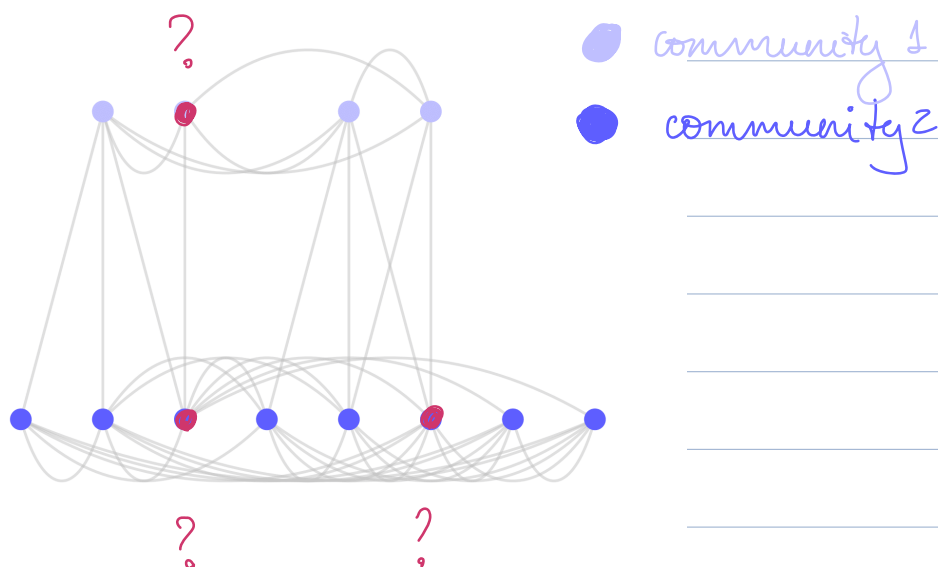
→ recap SBM

→ spectral clustering

→ spectral embeddings

→ contextual SBM

→ GNNs for community detection



► The stochastic block model (recap)

An n-node SBM graph with C communities is given by:

$$P = YBY^T$$

$$A \sim P \quad (\text{adjacency})$$

where: •) $Y \in \{0,1\}^{n \times C}$ is community assignment matrix; if $Y_{ic} = 1$, node i belongs to comm. c

•) B is the matrix of intra & inter-comm. probabilities; $B_{c_1 c_2}$ is the edge probability for a node pair (i,j) s.t. $Y_{ic_1} = 1, Y_{jc_2} = 1$

(or symmetric \rightarrow SSBM!)

→ We say the SBM is balanced when all comms. have the same size (n/C); it is unbalanced otherwise.

→ While B can take any value (as long as it's symmetric), often:

$$B_{c_1 c_2} = \begin{cases} p & \text{if } c_1 = c_2 \\ q & \text{if } c_1 \neq c_2 \end{cases}$$

- I.e., the probability of connecting to nodes in the same community is the same for all comms;
- the probability of connecting across communities is the same for all pairs c_1, c_2 w/ $c_1 \neq c_2$

→ Consider the 2-community balanced SBM w/ intra & intercommunity parameters p & q

$$B = \begin{bmatrix} p & q \\ q & p \end{bmatrix}$$

For an n -node graph, suppose the nodes are labeled such that the first $\frac{n}{2}$ are in community 1, and the remaining $\frac{n}{2}$ nodes are in community 2.

• Then,

$$Y = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \\ 0 & 1 \\ \vdots & \vdots \\ 0 & 1 \end{bmatrix} \begin{matrix} \frac{n}{2} \text{ community 1} \\ \frac{n}{2} \text{ community 2} \end{matrix}$$

and P is given by:

$$\left[\begin{array}{c|c} \underbrace{\begin{matrix} p & \dots & p \\ \vdots & & \vdots \\ p & \dots & p \\ q & \dots & q \\ \vdots & & \vdots \\ q & \dots & q \end{matrix}}_{n/2} & \underbrace{\begin{matrix} q & \dots & q \\ \vdots & & \vdots \\ q & \dots & q \\ p & \dots & p \\ \vdots & & \vdots \\ p & \dots & p \end{matrix}}_{n/2} \end{array} \right] \begin{matrix} \left. \vphantom{\begin{matrix} p & \dots & p \\ \vdots & & \vdots \\ p & \dots & p \\ q & \dots & q \\ \vdots & & \vdots \\ q & \dots & q \end{matrix}} \right\} n/2 \\ \left. \vphantom{\begin{matrix} q & \dots & q \\ \vdots & & \vdots \\ q & \dots & q \\ p & \dots & p \\ \vdots & & \vdots \\ p & \dots & p \end{matrix}} \right\} n/2 \end{matrix}$$

• Recall $A \sim P$, i.e.,

$$A_{ij} = \begin{cases} 1 & \text{with probability } p_{ij} \\ 0 & \text{with prob. } 1 - p_{ij} \end{cases}$$

• Hence, $E(A) = P$. Let's compute the eigenvalues & eigenvectors of $P/E(A)$

$$v_1 = \begin{bmatrix} 1 \\ \frac{1}{n} \\ \vdots \\ \frac{1}{n} \end{bmatrix} \quad \text{since the rows of } P \text{ all have the same sum}$$

2

$$\begin{pmatrix} \boxed{p} & \boxed{q} \\ \boxed{q} & \boxed{p} \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{n}} \\ \vdots \\ \frac{1}{\sqrt{n}} \end{pmatrix} = \frac{(p+q)}{2} \begin{pmatrix} \frac{1}{\sqrt{n}} \\ \vdots \\ \frac{1}{\sqrt{n}} \end{pmatrix}$$

$$\Rightarrow \lambda_1 = p+q$$

Further, note:

$$\begin{pmatrix} \boxed{p} & \boxed{q} \\ \boxed{q} & \boxed{p} \end{pmatrix} = \frac{(p+q)}{2} \begin{pmatrix} \overbrace{1, 1}^T \\ 1 \end{pmatrix} + \frac{(p-q)}{2} \begin{pmatrix} \boxed{1} & \boxed{-1} \\ \boxed{-1} & \boxed{1} \end{pmatrix}$$

$$= \frac{(p+q)}{2} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} [1 \dots 1]$$

$$+ \frac{(p-q)}{2} \begin{bmatrix} 1 \\ \vdots \\ 1 \\ -1 \\ \vdots \\ -1 \end{bmatrix} \begin{bmatrix} \overbrace{1 \dots 1}^{n/2} & \overbrace{-1 \dots -1}^{n/2} \end{bmatrix}$$

$$\Rightarrow v_2 = \left(\underbrace{\frac{1}{\sqrt{n}} \dots \frac{1}{\sqrt{n}}}_{n/2} \underbrace{-\frac{1}{\sqrt{n}} \dots -\frac{1}{\sqrt{n}}}_{n/2} \right)$$

$$\lambda_2 = \frac{p - q}{2}$$

The second eigenvector associated w/ the second largest eigenvalue (in abs. value) of $E(A) = P$

reveals the comm. assignment

This is the intuition behind spectral clustering;
 For an arbitrary graph with adjacency matrix A , we assume $A = A_{\text{SBM}} + \mathcal{E}$, where \mathcal{E} is random noise w/ $E(\mathcal{E}) = 0$; since $E(A) = E(A_{\text{SBM}}) = P$, we can estimate the community assignment from the eigenvectors of A , which is an unbiased estimator of P .

► Spectral clustering

↳ what happens for $C > 2$?

- the community assignment is a linear combination of the top C eigenvectors (in absolute value)

↳ suppose we are given A and C . How to estimate Y ?

The spectral clustering algorithm:

① Diagonalize $A = V \Lambda V^T$

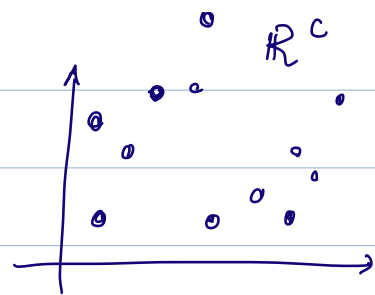
② Order eigenvectors by decreasing eigenvalue magnitude

$$V_C = \begin{bmatrix} | & | & | & & | \\ v_1 & v_2 & v_3 & \dots & v_C \\ | & | & | & & | \end{bmatrix} = \begin{bmatrix} \text{---} u_1 \text{---} \\ \text{---} u_2 \text{---} \\ \vdots \\ \text{---} u_n \text{---} \end{bmatrix} \quad \begin{matrix} u_j \in \mathbb{R}^C \\ v_i \in \mathbb{R}^n \end{matrix}$$

u_j is the embedding of node j in \mathbb{R}^C

③ Cluster the u_i into C groups.

There are multiple clustering algos;
most common is k-means clustering



► Semi-supervised learning with spectral embeddings

Spectral clustering is fully unsupervised - it only assumes knowledge of A .

In the case where the labels of a node subset $\mathcal{Y} \subset V$ are known, these can be used to enhance predictive power.

Explicitly, we aim to solve the problem:

(*)

$$\min_f \sum_{i \in \mathcal{Y}} \mathbb{I}([f(A)]_i = y_i)$$

\mathbb{I} indicator fn
 $[f(A)]_i$ some parametric function
 y_i one-hot comm. vector
 $[y_{ic}] = 1$ if i is in comm. c

(in practice the indicator/0-1 loss is substituted for the cross-entropy loss)

Spectral embeddings: $f(A) = \sigma(V_C W)$
(1-layer)

$W \in \mathbb{R}^{C \times C}$ linear map
 σ nonlinearity

↳ can be thought of as an FCNN on the C -dimensional embeddings u_1, \dots, u_n of nodes in \mathcal{V}

To find f^* , i.e., W^* , we solve (*) using gradient-descent methods

► Information-theoretic threshold

→ What happens when $p \approx q$ in the balanced SBM with $C=2$?

When $p = q$, we actually have an Erdős-Rényi graph, in which the edge probability is constant for all nodes; there are no communities to distinguish

But even when $p \neq q$, there is a region around $p = q$ where detection of communities is impossible in an information theoretic sense:

- Let $SNR = \frac{(p-q)^2}{2(p+q)}$ (signal-to-noise ratio)
degree discrepancy in \neq subgraphs \sim avg degree ("noise")

- If $SNR < 1/n$, almost exact recovery
$$\left\{ P \left(\frac{1}{n} \mathbb{I}([f(A)]_i = y_i) = 1 - o(1) \right) = 1 - o(1) \right\}$$

is impossible. Even with infinite time & resources, there is no algorithm that can recover the true communities given A .

For the pf, check the works of Massoulié (2014)
Mossel (2014)
Abbe (survey)

E.g.: sparse graphs

$$p = \frac{a}{n} ; \quad q = \frac{b}{n}$$

$$\frac{|E|}{n^2} \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

(avg. degree vanishes)

$$SNR = \frac{\left(\frac{a}{n} - \frac{b}{n}\right)^2}{2\left(\frac{a+b}{n}\right)} = \frac{1}{2n} \frac{(a-b)^2}{(a+b)}$$