

Community recovery: Recap

Definition 1. Let n be a positive integer (the number of vertices), k be a positive integer (the number of communities), $p = (p_1, \dots, p_k)$ be a probability vector on $[k] := \{1, \dots, k\}$ (the prior on the k communities) and W be a $k \times k$ symmetric matrix with entries in $[0, 1]$ (the connectivity probabilities). The pair (X, G) is drawn under $\text{SBM}(n, p, W)$ if X is an n -dimensional random vector with i.i.d. components distributed under p , and G is an n -vertex simple graph where vertices i and j are connected with probability W_{X_i, X_j} , independently of other pairs of vertices. We also define the community sets by $\Omega_i = \Omega_i(X) := \{v \in [n] : X_v = i\}$, $i \in [k]$.

Definition 2. (X, G) is drawn under $\text{SSBM}(n, k, q_{\text{in}}, q_{\text{out}})$ if W takes value q_{in} on the diagonal and q_{out} off the diagonal, and if the community prior is $p = \{1/k\}^k$ in the Bernoulli model, and X is drawn uniformly at random with the constraints $|\{v \in [n] : X_v = i\}| = n/k$, n a multiple of k , in the uniform or strictly balanced model.

Definition 3 (Agreement and normalized agreement). *The agreement between two community vectors $x, y \in [k]^n$ is obtained by maximizing the common components between x and any relabelling of y , i.e.,*

$$A(x, y) = \max_{\pi \in S_k} \frac{1}{n} \sum_{i=1}^n \mathbb{1}(x_i = \pi(y_i)), \quad (9)$$

Definition 4. *Let $(X, G) \sim \text{SBM}(n, p, W)$. The following recovery requirements are solved if there exists an algorithm that takes G as an input and outputs $\hat{X} = \hat{X}(G)$ such that*

- **Exact recovery:** $\mathbb{P}\{A(X, \hat{X}) = 1\} = 1 - o(1)$,
- **Almost exact recovery:** $\mathbb{P}\{A(X, \hat{X}) = 1 - o(1)\} = 1 - o(1)$,

Theorem 3. *[ABH14, MNS14a] Exact recovery in $\text{SSBM}(n, 2, a \log(n)/n, b \log(n)/n)$ is solvable and efficiently so if $|\sqrt{a} - \sqrt{b}| > \sqrt{2}$ and unsolvable if $|\sqrt{a} - \sqrt{b}| < \sqrt{2}$.*

Theorem 3. *[ABH14, MNS14a] Exact recovery in $\text{SSBM}(n, 2, a \log(n)/n, b \log(n)/n)$ is solvable and efficiently so if $|\sqrt{a} - \sqrt{b}| > \sqrt{2}$ and unsolvable if $|\sqrt{a} - \sqrt{b}| < \sqrt{2}$.*

In the two-community case, denoting by N_{in} and N_{out} the number of edges inside and across the clusters respectively,

$$\mathbb{P}\{G = g|\Omega = \omega\} \propto \left(\frac{q_{out}(1 - q_{in})}{q_{in}(1 - q_{out})}\right)^{N_{out}}. \quad (19)$$

Assuming $q_{in} \geq q_{out}$, we have $\frac{q_{out}(1-q_{in})}{q_{in}(1-q_{out})} \leq 1$ and thus

MAP is equivalent to finding a min-bisection of G ,

i.e., a balanced partition with the least number of crossing edges.

Spectral relaxations. Consider again the symmetric SBM with strictly balanced communities. Recall that MAP maximizes

$$\max_{\substack{x \in \{+1, -1\}^n \\ x^t \mathbf{1}^n = 0}} x^t A x, \tag{20}$$

since this counts the number of edges inside the clusters minus the number of edges across the clusters, which is equivalent to the min-bisection problem (the total number of edges being fixed by the graph). The general idea behind spectral methods is to

relax the integral constraint to an Euclidean constraint on real valued vectors. This leads to looking for a maximizer of

$$\max_{\substack{x \in \mathbb{R}^n: \|x\|_2^2 = n \\ x^t \mathbf{1}^n = 0}} x^t A x. \quad (21)$$

Without the constraint $x^t \mathbf{1}^n = 0$, the above maximization gives precisely the eigenvector corresponding to the largest eigenvalue of A . Note that $A \mathbf{1}^n$ is the vector containing the degrees of each node in g , and when g is an instance of the symmetric SBM, this concentrates to the same value for each vertex, and $\mathbf{1}^n$ is close to an eigenvector of A . Since A is real symmetric, this suggests that the constraint $x^t \mathbf{1}^n = 0$ leads the maximization (21) to focus on the eigenspace orthogonal to the first eigenvector, and thus to the eigenvector corresponding to the second largest eigenvalue. Thus it is reasonable to take the second largest eigenvector $\phi_2(A)$ of A and round it to obtain an efficient relaxation of MAP:

$$\hat{X}_{\text{spec}} = \begin{cases} 1 & \text{if } \phi_2(A) \geq 0, \\ 2 & \text{if } \phi_2(A) < 0. \end{cases} \quad (22)$$

Now we describe a spectral method. To simplify presentation, it is assumed without loss of generality that: $x_i^* = 1$ for any $1 \leq i \leq n/2$, and $x_i^* = -1$ for any $i > n/2$.

A starting point for the algorithm design is to examine the mean of the adjacency matrix, given as follows

$$\mathbb{E}[\mathbf{A}] = \begin{bmatrix} p \mathbf{1}_{n/2} \mathbf{1}_{n/2}^\top & q \mathbf{1}_{n/2} \mathbf{1}_{n/2}^\top \\ q \mathbf{1}_{n/2} \mathbf{1}_{n/2}^\top & p \mathbf{1}_{n/2} \mathbf{1}_{n/2}^\top \end{bmatrix} - p \mathbf{I}.$$

As revealed by the above calculation, the matrix constructed below

$$\mathbf{M} = \mathbf{A} - \frac{p+q}{2} \mathbf{1}_n \mathbf{1}_n^\top + p \mathbf{I} \tag{3.30}$$

exhibits an approximate rank-1 structure, in the sense that its mean

$$\mathbf{M}^* := \mathbb{E}[\mathbf{M}] = \frac{p-q}{2} \begin{bmatrix} \mathbf{1}_{n/2} \\ -\mathbf{1}_{n/2} \end{bmatrix} \begin{bmatrix} \mathbf{1}_{n/2}^\top & -\mathbf{1}_{n/2}^\top \end{bmatrix} \tag{3.31}$$

is a rank-1 matrix. The leading eigenvalue of \mathbf{M}^* and its associated eigenvector are given respectively by

$$\lambda^* := \frac{(p-q)n}{2}, \quad \text{and} \quad \mathbf{u}^* := \frac{1}{\sqrt{n}} \begin{bmatrix} \mathbf{1}_{n/2} \\ -\mathbf{1}_{n/2} \end{bmatrix}. \quad (3.32)$$

Crucially, the eigenvector \mathbf{u}^* encapsulates the precise community structure we seek to recover: all positive entries of \mathbf{u}^* correspond to vertices from one community, while the remaining ones form another community.

Inspired by the above calculation, a candidate spectral clustering algorithm consists of eigendecomposition followed by entrywise rounding:

1. Compute the leading eigenvector \mathbf{u} of \mathbf{M} (constructed in (3.30));
2. Compute the estimate $\mathbf{x} = [x_i]_{1 \leq i \leq n}$ such that for any $1 \leq i \leq n$,

$$x_i = \text{sgn}(u_i) = \begin{cases} 1, & \text{if } u_i > 0, \\ -1, & \text{if } u_i \leq 0. \end{cases} \quad (3.33)$$

Theorem 3.8. Consider the setting in Section 3.4.1, and suppose that

$$p \gtrsim \frac{\log n}{n}, \quad \text{and} \quad \sqrt{\frac{p}{n}} = o(p - q). \quad (3.35)$$

With probability exceeding $1 - O(n^{-8})$, the spectral method achieves

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i = x_i^*\} = 1 - o(1), \quad \text{or} \quad \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i = -x_i^*\} = 1 - o(1).$$

Special case 1: $p - q \gg \sqrt{\log n} / n$, if $p \asymp (\log n) / n$.

Special case 2: $p - q \gg 1/\sqrt{n}$, if $p \asymp 1$,

Theorem 3.4. Consider a *symmetric* random matrix $\mathbf{X} = [X_{i,j}]_{1 \leq i,j \leq n}$ in $\mathbb{R}^{n \times n}$, whose entries are independently generated and obey

$$\mathbb{E}[X_{i,j}] = 0, \quad \text{and} \quad |X_{i,j}| \leq B, \quad 1 \leq i, j \leq n. \quad (3.6)$$

Define

$$\nu := \max_i \sum_j \mathbb{E}[X_{i,j}^2]. \quad (3.7)$$

Then there exists some universal constant $c > 0$ such that for any $t \geq 0$,

$$\mathbb{P}\left\{\|\mathbf{X}\| \geq 4\sqrt{\nu} + t\right\} \leq n \exp\left(-\frac{t^2}{cB^2}\right). \quad (3.8)$$

This result, which appeared in Bandeira and Van Handel (2016, Remark 3.13), can be established

- As a useful corollary, if we know *a priori* that $\mathbb{E}[X_{i,j}^2] \leq \sigma^2$ for all $1 \leq i, j \leq n$, then Theorem 3.4 implies that

$$\|\mathbf{X}\| \leq 4\sigma\sqrt{n} + \tilde{c}B\sqrt{\log n} \quad (3.9)$$

with probability at least $1 - n^{-8}$ for some constant $\tilde{c} > 0$. To see this, it suffices to set $\tilde{c} = \sqrt{9c}$ and take $t = B\sqrt{9c \log n}$ in (3.8).

2.3.1 Setup and notation

Let M^* and $M = M^* + E$ be two $n \times n$ real symmetric matrices. We express the eigendecomposition of M^* and M as follows

$$M^* = \sum_{i=1}^n \lambda_i^* \mathbf{u}_i^* \mathbf{u}_i^{*\top} = \begin{bmatrix} U^* & U_{\perp}^* \end{bmatrix} \begin{bmatrix} \Lambda^* & \mathbf{0} \\ \mathbf{0} & \Lambda_{\perp}^* \end{bmatrix} \begin{bmatrix} U^{*\top} \\ U_{\perp}^{*\top} \end{bmatrix}; \quad (2.11)$$

$$M = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^{\top} = \begin{bmatrix} U & U_{\perp} \end{bmatrix} \begin{bmatrix} \Lambda & \mathbf{0} \\ \mathbf{0} & \Lambda_{\perp} \end{bmatrix} \begin{bmatrix} U^{\top} \\ U_{\perp}^{\top} \end{bmatrix}. \quad (2.12)$$

Here, $\{\lambda_i\}$ (resp. $\{\lambda_i^*\}$) denote the eigenvalues of M (resp. M^*), and \mathbf{u}_i (resp. \mathbf{u}_i^*) stands for the eigenvector associated with the eigenvalue λ_i (resp. λ_i^*). Additionally, we take

$$\begin{aligned} U &:= [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{R}^{n \times r}, & U_{\perp} &:= [\mathbf{u}_{r+1}, \dots, \mathbf{u}_n] \in \mathbb{R}^{n \times (n-r)}, \\ \Lambda &:= \text{diag}([\lambda_1, \dots, \lambda_r]), & \Lambda_{\perp} &:= \text{diag}([\lambda_{r+1}, \dots, \lambda_n]). \end{aligned}$$

The matrices U^* , U_{\perp}^* , Λ^* , and Λ_{\perp}^* are defined analogously.

$$\text{dist}(U, U^*) := \min_{R \in \mathcal{O}^{r \times r}} \|UR - U^*\|; \quad (2.10a)$$

$$\text{dist}_{\mathbb{F}}(U, U^*) := \min_{R \in \mathcal{O}^{r \times r}} \|UR - U^*\|_{\mathbb{F}}. \quad (2.10b)$$

Corollary 2.8. Consider the settings in Section 2.3.1. Suppose that $|\lambda_1^*| \geq |\lambda_2^*| \geq \cdots \geq |\lambda_r^*| > |\lambda_{r+1}^*| \geq \cdots \geq |\lambda_n^*|$ and $|\lambda_1| \geq |\lambda_2| \geq \cdots \geq |\lambda_n|$ (i.e., the eigenvalues are sorted by their magnitudes). If $\|\mathbf{E}\| < (1 - 1/\sqrt{2})(|\lambda_r^*| - |\lambda_{r+1}^*|)$, then

$$\text{dist}(\mathbf{U}, \mathbf{U}^*) \leq \sqrt{2} \|\sin \Theta\| \leq \frac{2\|\mathbf{E}\mathbf{U}^*\|}{|\lambda_r^*| - |\lambda_{r+1}^*|} \leq \frac{2\|\mathbf{E}\|}{|\lambda_r^*| - |\lambda_{r+1}^*|}; \quad (2.18a)$$

$$\text{dist}_F(\mathbf{U}, \mathbf{U}^*) \leq \sqrt{2} \|\sin \Theta\|_F \leq \frac{2\|\mathbf{E}\mathbf{U}^*\|_F}{|\lambda_r^*| - |\lambda_{r+1}^*|} \leq \frac{2\sqrt{r}\|\mathbf{E}\|}{|\lambda_r^*| - |\lambda_{r+1}^*|}. \quad (2.18b)$$

Theorem 4.6. Fix any constant $\varepsilon > 0$, and consider the setting of Section 3.4.1. Suppose $p = \frac{\alpha \log n}{n}$ and $q = \frac{\beta \log n}{n}$ for some sufficiently large constants $\alpha > \beta > 0$.¹ In addition, assume that

$$(\sqrt{p} - \sqrt{q})^2 \geq 2(1 + \varepsilon) \frac{\log n}{n}. \quad (4.45)$$

With probability $1 - o(1)$, the spectral method in Section 3.4.2 yields

$$x_i = x_i^* \text{ for all } 1 \leq i \leq n, \quad \text{or} \quad x_i = -x_i^* \text{ for all } 1 \leq i \leq n.$$

Next, we turn to vector and matrix norms. For any vector \mathbf{v} , we denote by $\|\mathbf{v}\|_2$, $\|\mathbf{v}\|_1$ and $\|\mathbf{v}\|_\infty$ its ℓ_2 norm, ℓ_1 norm and ℓ_∞ norm, respectively. For any matrix $\mathbf{A} = [A_{i,j}]_{1 \leq i \leq m, 1 \leq j \leq n}$, we let $\|\mathbf{A}\|$, $\|\mathbf{A}\|_*$, $\|\mathbf{A}\|_F$ and $\|\mathbf{A}\|_\infty$ represent respectively its spectral norm (i.e., the largest singular value of \mathbf{A}), its nuclear norm (i.e., the sum of singular values of \mathbf{A}), its Frobenius norm (i.e., $\|\mathbf{A}\|_F := \sqrt{\sum_{i,j} A_{i,j}^2}$), and its entrywise ℓ_∞ norm (i.e., $\|\mathbf{A}\|_\infty := \max_{i,j} |A_{i,j}|$). We also refer to $\|\mathbf{A}\|_{2,\infty}$ as the $\ell_{2,\infty}$ norm of \mathbf{A} , defined as $\|\mathbf{A}\|_{2,\infty} := \max_i \|\mathbf{A}_{i,\cdot}\|_2$. Similarly, we define the $\ell_{\infty,2}$ norm of \mathbf{A} as $\|\mathbf{A}\|_{\infty,2} := \|\mathbf{A}^\top\|_{2,\infty}$. In addition, for any matrices $\mathbf{A} = [A_{i,j}]_{1 \leq i \leq m, 1 \leq j \leq n}$ and $\mathbf{B} = [B_{i,j}]_{1 \leq i \leq m, 1 \leq j \leq n}$, the inner product of \mathbf{A} and \mathbf{B} is defined as and denoted by $\langle \mathbf{A}, \mathbf{B} \rangle = \sum_{1 \leq i \leq m, 1 \leq j \leq n} A_{i,j} B_{i,j} = \text{Tr}(\mathbf{A}^\top \mathbf{B})$.

Ground truth. Consider a rank- r symmetric matrix $\mathbf{M}^* \in \mathbb{R}^{n \times n}$ with eigenvectors $\mathbf{u}_1^*, \dots, \mathbf{u}_n^*$ and associated eigenvalues $\lambda_1^*, \dots, \lambda_n^*$ obeying

$$|\lambda_1^*| \geq |\lambda_2^*| \geq \dots \geq |\lambda_r^*| > 0 \quad \text{and} \quad \lambda_{r+1}^* = \dots = \lambda_n^* = 0. \quad (4.22)$$

We shall write the eigendecomposition $\mathbf{M}^* = \mathbf{U}^* \mathbf{\Lambda}^* \mathbf{U}^{*\top}$ as usual, where $\mathbf{\Lambda}^* := \text{diag}([\lambda_1^*, \dots, \lambda_r^*])$ and $\mathbf{U}^* := [\mathbf{u}_1^*, \dots, \mathbf{u}_r^*] \in \mathbb{R}^{n \times r}$. Denote the condition number of \mathbf{M}^* as

$$\kappa := |\lambda_1^*| / |\lambda_r^*|. \quad (4.23)$$

Akin to Definition 3.1, the incoherence parameter of \mathbf{M}^* is defined as

$$\mu := \frac{n \|\mathbf{U}^*\|_{2,\infty}^2}{r}, \quad (4.24)$$

a parameter that captures how well the energy of \mathbf{U}^* is spread out across all rows and that obeys (see Remark 3.12)

$$1 \leq \mu \leq n/r. \quad (4.25)$$

Observed data. What we observe is a corrupted version

$$\mathbf{M} = \mathbf{M}^* + \mathbf{E} \in \mathbb{R}^{n \times n}, \quad (4.26)$$

where \mathbf{E} is a symmetric noise matrix. We denote by $\{\lambda_i\}_{1 \leq i \leq n}$ the set of eigenvalues of \mathbf{M} obeying

$$|\lambda_1| \geq |\lambda_2| \geq \cdots \geq |\lambda_n|, \quad (4.27)$$

and let \mathbf{u}_i be the eigenvector of \mathbf{M} associated with λ_i . We shall introduce the diagonal matrix $\mathbf{\Lambda} \in \mathbb{R}^{r \times r}$ as $\mathbf{\Lambda} := \text{diag}([\lambda_1, \cdots, \lambda_r])$.

Noise assumptions. This section aims to cover a fairly broad class of scenarios of independent noise. In particular, the noise matrix considered herein is assumed to satisfy the mild conditions listed below.

Assumption 4.1. The entries in the lower triangular part of $\mathbf{E} = [E_{i,j}]_{1 \leq i,j \leq n}$ are independently generated obeying

$$\mathbb{E}[E_{i,j}] = 0, \quad \mathbb{E}[E_{i,j}^2] =: \sigma_{i,j}^2 \leq \sigma^2, \quad |E_{i,j}| \leq B, \quad \text{for all } i \geq j. \quad (4.28)$$

In particular, σ^2 is taken to be the smallest choice satisfying (4.28). Further, it is assumed that

$$c_b := \frac{B}{\sigma \sqrt{n}/(\mu \log n)} = O(1). \quad (4.29)$$

Goal and algorithm. We seek to estimate \mathbf{U}^* based on \mathbf{M} . Towards this, a simple spectral method computes the matrix $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{R}^{n \times r}$ that comprises the top- r leading eigenvectors of \mathbf{M} .

4.2.2 $\ell_{2,\infty}$ and ℓ_∞ theoretical guarantees

The leave-one-out argument introduced before, when properly strengthened, enables powerful $\ell_{2,\infty}$ performance guarantees for the spectral estimate \mathbf{U} , which concern row-wise perturbation of the eigenspace. Before continuing, we remind the readers of the global rotation ambiguity, namely, in general we cannot expect \mathbf{U} to be close to \mathbf{U}^* unless suitable global rotation is taken into account. In light of this, we introduce the following notation that helps identify a proper rotation matrix.

Definition 4.1. For any matrix \mathbf{Z} with SVD $\mathbf{Z} = \mathbf{U}_Z \boldsymbol{\Sigma}_Z \mathbf{V}_Z^\top$ (where \mathbf{U}_Z and \mathbf{V}_Z represent respectively the left and right singular matrices of \mathbf{Z} , and $\boldsymbol{\Sigma}_Z$ is a diagonal matrix composed of the singular values), define

$$\text{sgn}(\mathbf{Z}) := \mathbf{U}_Z \mathbf{V}_Z^\top \tag{4.30}$$

to be the matrix sign function of \mathbf{Z} .

Remark 4.1. The matrix sign function is commonly encountered when aligning two matrices—classically known as the orthogonal Procrustes problem (Schönemann, 1966). Consider any two matrices $\widehat{\mathbf{B}}, \mathbf{B} \in \mathbb{R}^{n \times r}$ with $r \leq n$. Among all rotation matrices, the one that best aligns $\widehat{\mathbf{B}}$ with \mathbf{B} is precisely $\text{sgn}(\widehat{\mathbf{B}}^\top \mathbf{B})$ (see, e.g., (Ma *et al.*, 2020, Appendix D.2.1)), namely,

$$\text{sgn}(\widehat{\mathbf{B}}^\top \mathbf{B}) = \arg \min_{\mathbf{O} \in \mathcal{O}^{r \times r}} \|\widehat{\mathbf{B}} \mathbf{O} - \mathbf{B}\|_{\text{F}}^2.$$

Theorem 4.2. Consider the settings and assumptions in Section 4.2.1. Define $\mathbf{H} := \mathbf{U}^\top \mathbf{U}^*$. With probability exceeding $1 - O(n^{-5})$, one has

$$\|\mathbf{U} \operatorname{sgn}(\mathbf{H}) - \mathbf{U}^*\|_{2,\infty} \lesssim \frac{\sigma \kappa \sqrt{\mu r} + \sigma \sqrt{r \log n}}{|\lambda_r^*|}, \quad (4.31a)$$

$$\begin{aligned} & \|\mathbf{U} \operatorname{sgn}(\mathbf{H}) - \mathbf{M} \mathbf{U}^* (\mathbf{\Lambda}^*)^{-1}\|_{2,\infty} \\ & \lesssim \frac{\sigma \kappa \sqrt{\mu r}}{|\lambda_r^*|} + \frac{\sigma^2 \sqrt{r n \log n} + \sigma B \sqrt{\mu r \log^3 n}}{(\lambda_r^*)^2}, \end{aligned} \quad (4.31b)$$

provided that $\sigma \sqrt{n \log n} \leq c_\sigma |\lambda_r^*|$ for some sufficiently small constant $c_\sigma > 0$.

