

Graph Clustering (GC)

- ▶ **Graph Clustering (GC)** is a core analysis technique used for network data:
 - ▶ Social Networks, Ecological Networks, Transportation Networks, Brain Networks etc.
- ▶ Real networks are often available with **partial observation of its edges** due to:
 - ▶ Massive Data, Cost, Security/Privacy



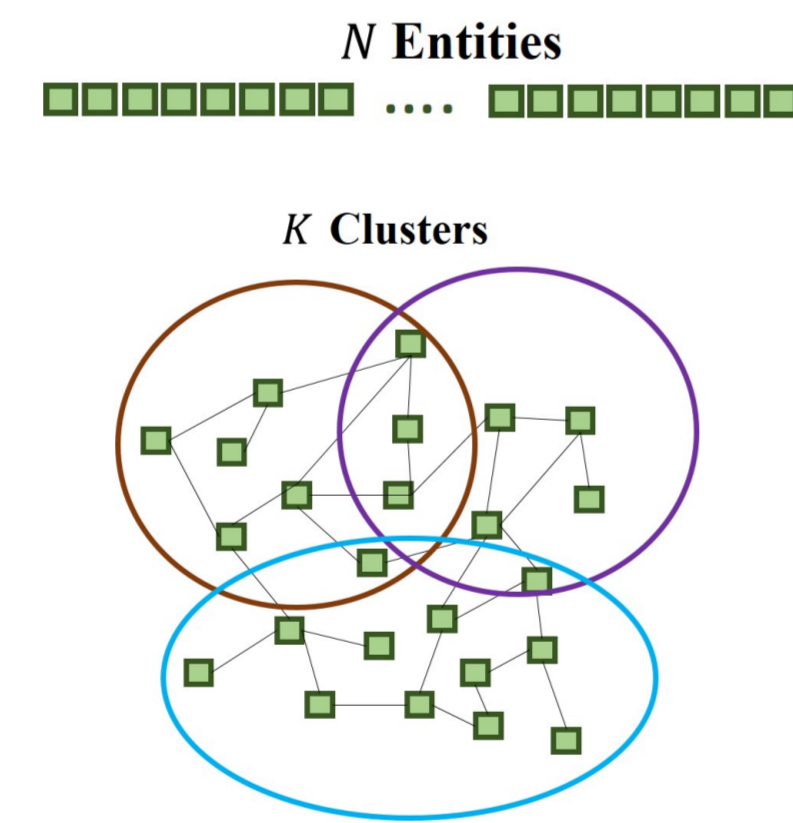
Existing Work with Provable Guarantees

- A number of works [Korlakai Vinayak et al., 2014; Korlakai Vinayak and Hassibi, 2016; Chen et al., 2014], which proposed GC under partial edge observation, features
- **single membership identification**
 - ▶ the entities often admit mixed membership in real-world networks
 - **random query based edge acquisition scheme**
 - ▶ may not be easy to implement in some applications; e.g., in field surveys and in networks with hidden or intentionally removed edges
 - **convex optimization based problem formulation**
 - ▶ hard to scale up for real-world large graphs

We aim to design a **systematic edge query scheme** for mixed membership identification via a **lightweight algorithm** with **provable guarantees**.

Mixed Membership Model

- ▶ The n th entity belongs to k th cluster with prob. m_{kn}
 - ▶ $\sum_{k=1}^K m_{k,n} = 1$, $m_{k,n} \geq 0$.
- ▶ $\mathbf{m}_n = [m_{1,n}, \dots, m_{K,n}]^\top$ is called as the **membership vector** of n .
- ▶ $\mathbf{M} = [\mathbf{m}_1, \dots, \mathbf{m}_N] \in \mathbb{R}^{K \times N}$ is called as the **membership matrix**.
- ▶ $\mathbf{B} \in \mathbb{R}^{K \times K}$ is **cluster-cluster interaction matrix**.

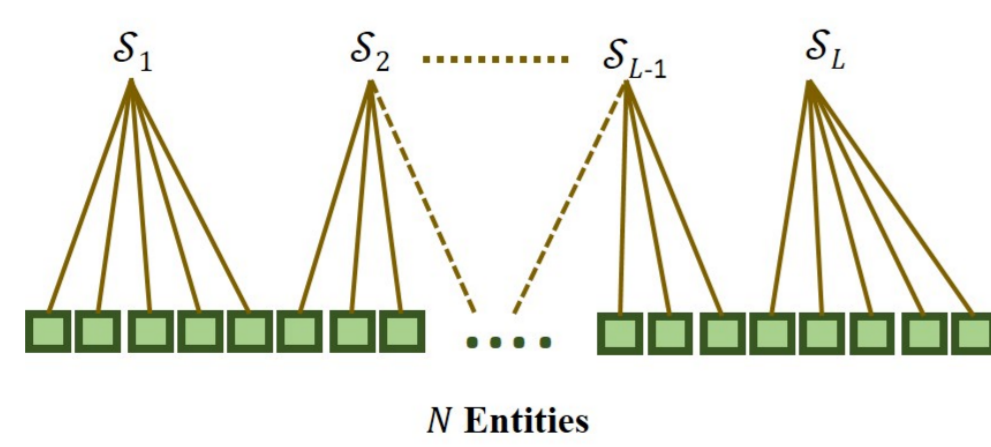


- ▶ The edges of the graph are represented using **adjacency matrix** $\mathbf{A} \in \{0, 1\}^{N \times N}$:
$$\mathbf{A}(i, j) \sim \text{Bernoulli}(\mathbf{P}(i, j)), \quad \mathbf{P} = \mathbf{M}^\top \mathbf{B} \mathbf{M}, \quad \mathbf{1}^\top \mathbf{M} = \mathbf{1}^\top, \quad \mathbf{M} \geq 0.$$

Proposed Systematic Edge Query

$$\mathcal{S}_1 \cup \dots \cup \mathcal{S}_L = \{1, \dots, N\}$$

$$\mathcal{S}_\ell \cap \mathcal{S}_m = \emptyset, \quad \forall \ell \neq m$$



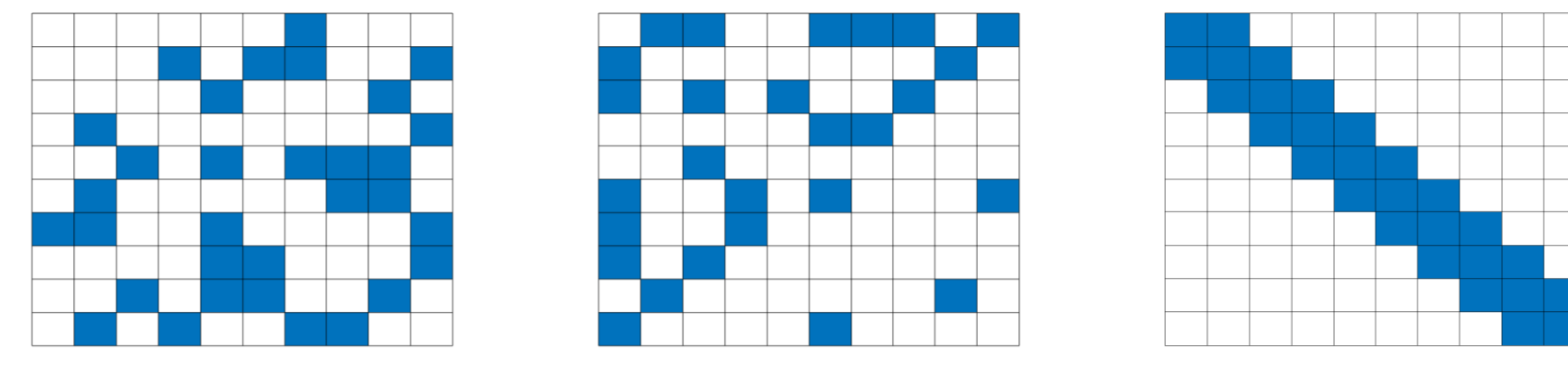
Adjacency Submatrix between \mathcal{S}_ℓ and $\mathcal{S}_m \implies \mathbf{A}_{\ell,m} \in \mathbb{R}^{|\mathcal{S}_\ell| \times |\mathcal{S}_m|}$

Edge Query Principle (EQP)

- For every $\ell \in [L]$, $K \leq |\mathcal{S}_\ell|$ holds. Let $m_r \in [L]$ and $\{\ell_r\}_{r=1}^L = [L]$.
- For every ℓ_r , there exists a pair of indices m_r and ℓ_{r+1} where $\ell_{r+1} \neq \ell_r$ such that the edges from the blocks \mathbf{A}_{ℓ_r, m_r} and $\mathbf{A}_{\ell_{r+1}, m_r}$ are queried.

Algorithm Design for Learning M under EQP

Some examples for EQP patterns with $N = 1000$, $K = 5$ and $L = 10$.

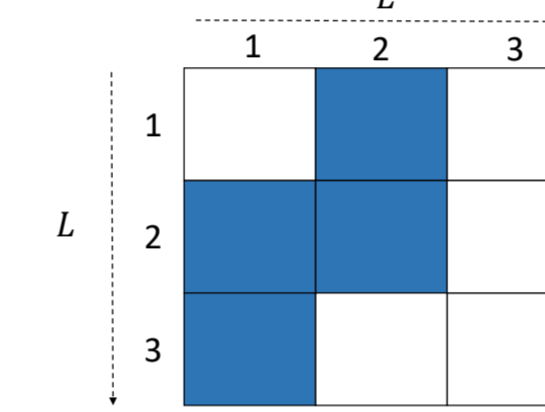


- Step 1:** Estimate $\mathbf{U} \in \mathbb{R}^{N \times K}$ such that $\text{range}(\mathbf{U}) = \text{range}(\mathbf{M}^\top)$
- Step 2:** Estimate \mathbf{M} from \mathbf{U} via structured matrix factorization (SMF)

Consider $L = 3$ and $\mathbf{A}_{\ell,m} = \mathbf{P}_{\ell,m} = \mathbf{M}_\ell^\top \mathbf{B} \mathbf{M}_m$:

$$\mathbf{P}_{1,2} = \mathbf{M}_1^\top \mathbf{B} \mathbf{M}_2, \quad \mathbf{P}_{2,2} = \mathbf{M}_2^\top \mathbf{B} \mathbf{M}_2,$$

$$\mathbf{P}_{2,1} = \mathbf{M}_2^\top \mathbf{B} \mathbf{M}_1, \quad \mathbf{P}_{3,1} = \mathbf{M}_3^\top \mathbf{B} \mathbf{M}_1.$$



- ▶ Define $\mathbf{C}_1 := [\mathbf{P}_{1,2}, \mathbf{P}_{2,2}]^\top$ and $\mathbf{C}_2 := [\mathbf{P}_{2,1}, \mathbf{P}_{3,1}]^\top$. Consider their top- K SVD:

$$\mathbf{C}_1 = [\mathbf{U}_1^\top, \mathbf{U}_2^\top]^\top \Sigma \mathbf{V}^\top, \quad \mathbf{C}_2 = [\mathbf{U}_2^\top, \mathbf{U}_3^\top]^\top \Sigma \mathbf{V}^\top.$$

- ▶ The bases of $\text{range}(\mathbf{M}_1^\top)$, $\text{range}(\mathbf{M}_2^\top)$ and $\text{range}(\mathbf{M}_3^\top)$ are:
$$\mathbf{U}_1 = \mathbf{M}_1^\top \mathbf{B} \Theta, \quad \mathbf{U}_2 = \mathbf{M}_2^\top \mathbf{B} \Theta, \quad \mathbf{U}_3 = \mathbf{M}_3^\top \mathbf{B} \Phi, \quad \Phi \neq \Theta \text{ in general.}$$
- ▶ We need a certain \mathbf{U}_3 such that $\text{range}([\mathbf{U}_1^\top, \mathbf{U}_2^\top, \mathbf{U}_3^\top]^\top) = \text{range}([\mathbf{M}_1, \mathbf{M}_2, \mathbf{M}_3]^\top)$.
- ▶ We can obtain such \mathbf{U}_3 as below:

$$\mathbf{U}_3 := \mathbf{U}_3 \mathbf{U}_2^\dagger \mathbf{U}_2 = \mathbf{M}_3^\top \mathbf{B} \Phi \times (\mathbf{M}_2^\top \mathbf{B} \Phi)^\dagger \times \mathbf{M}_2^\top \mathbf{B} \Theta = \mathbf{M}_3^\top \mathbf{B} \Theta.$$

- ▶ To estimate \mathbf{U} , this "subspace stitching" idea is recursively applied over the queried blocks \mathbf{A}_{ℓ_r, m_r} and $\mathbf{A}_{\ell_{r+1}, m_r}$ for $r = 1, \dots, L-1$.

$$\mathbf{U}^\top = \mathbf{G} \mathbf{M}, \quad \mathbf{M} \geq 0, \quad \mathbf{1}^\top \mathbf{M} = \mathbf{1}^\top, \quad \mathbf{G} \in \mathbb{R}^{K \times K} \text{ is nonsingular.}$$

- ▶ **Successive Projection Algorithm (SPA)** [Gillis and Vavasis, 2014] can provably identify \mathbf{M} in K steps, if \mathbf{G} is nonsingular and if there exists $\{n_1, \dots, n_K\}$ such that $\mathbf{M}(:, n_k) = \mathbf{e}_k$ (**pure nodes**).

Algorithm 1: Proposed Algorithm

```

input : {A_{m, \ell}}, L, K
1 divide the blocks as {A_{\ell_r, m_r}}_{r=1}^L, {A_{\ell_{r+1}, m_r}}_{r=1}^{L-1}
   (where \ell_r \neq \ell_{r+1}, {\ell_r}_{r=1}^L = [L], m_r \in [L])
2 T \leftarrow \lfloor L/2 \rfloor;
3 C_T \leftarrow [A_{\ell_r, m_r}^\top, A_{\ell_{r+1}, m_r}^\top]^\top;
4 [\tilde{U}_r^\top, \tilde{U}_{r+1}^\top]^\top \Sigma_r \mathbf{V}_{m_r}^\top \leftarrow \text{svd}_K(C_T);
5 U_{ref} \leftarrow U_{\ell_{r+1}};
6 for r = T + 1 : 1 : L - 1 do
7   C_r \leftarrow [A_{\ell_r, m_r}^\top, A_{\ell_{r+1}, m_r}^\top]^\top;
8   [\tilde{U}_r^\top, \tilde{U}_{r+1}^\top]^\top \Sigma_r \mathbf{V}_{m_r}^\top \leftarrow \text{svd}_K(C_r);
9   U_{\ell_{r+1}} \leftarrow \tilde{U}_{\ell_{r+1}} \tilde{U}_r^\dagger U_{ref};
10  U_{ref} \leftarrow U_{\ell_{r+1}};
11 end
12 U_{ref} \leftarrow U_{\ell_T};
13 for r = T : -1 : 2 do
14  C_r \leftarrow [A_{\ell_r, m_r}^\top, A_{\ell_{r-1}, m_r}^\top]^\top;
15  [\tilde{U}_r^\top, \tilde{U}_{r-1}^\top]^\top \Sigma_r \mathbf{V}_{m_r}^\top \leftarrow \text{svd}_K(C_r);
16  U_{\ell_{r-1}} \leftarrow \tilde{U}_{\ell_{r-1}} \tilde{U}_r^\dagger U_{ref};
17  U_{ref} \leftarrow U_{\ell_{r-1}};
18 end
19 \hat{U} \leftarrow [U_1^\top, \dots, U_L^\top]^\top;
20 apply SPA on \hat{U} to estimate \hat{M}.
output: Estimated membership matrix \hat{M}.

```

Identifiability Results

Proposition 1: (Subspace Identifiability - Ideal Case)

Assume that $\mathbf{A}_{\ell,m} = \mathbf{P}_{\ell,m} = \mathbf{M}_\ell^\top \mathbf{B} \mathbf{M}_m \in \mathbb{R}^{|\mathcal{S}_\ell| \times |\mathcal{S}_m|}$ holds true for all $\ell, m \in [L]$ and $\text{rank}(\mathbf{M}) = \text{rank}(\mathbf{B}) = K$. Suppose that the $\mathbf{A}_{\ell,m}$'s are queried according to the proposed EQP. Then, the output $\hat{\mathbf{U}}$ by Algorithm 1 satisfies $\text{range}(\hat{\mathbf{U}}) = \text{range}(\mathbf{M}^\top)$.

Proposition 2: (Subspace Identifiability - Binary Observation Case)

Let $\rho := \max_{i,j} \mathbf{P}(i, j)$ be the maximal entry of \mathbf{P} . Suppose that $\rho = \Omega(L \log(N/L)/N)$ and $L = O(\rho N/d)$ where d is the maximal degree of all the nodes. Also assume that $N = \Omega\left(\max\left(L^2, \frac{(K\gamma^2)^L \rho \kappa^2(\mathbf{B})}{\sigma_{\min}^2(\mathbf{B})}\right)\right)$. Then, the output $\hat{\mathbf{U}}$ satisfies the following with probability of at least $1 - O(L^2/N)$:

$$\|\hat{\mathbf{U}} - \mathbf{U} \mathbf{O}\|_F = O\left(\frac{(K\gamma^2)^{L/2} \kappa(\mathbf{B}) \sqrt{\rho}}{\sigma_{\min}(\mathbf{B}) \sqrt{N/L}}\right),$$

where $\mathbf{O} \in \mathbb{R}^{K \times K}$ is an orthogonal matrix.

Larger L makes the error bound looser, but larger L means that only fewer queries need to be made, and thus less resource consuming.

Experiment Results

Synthetic Data Experiment:

- ▶ Baselines: GeoNMF [Mao et al., 2017], CD-MVSI [Mao et al., 2017]
- ▶ Parameters: $K = 5$, $L = 10$ with diagonal query pattern

Graph Size	Ideal Case ($\mathbf{A} = \mathbf{P}$)		Binary Observation Case		
	Proposed	Dist	Proposed	GeoNMF	CD-MVSI
N			Dist	MSE	MSE
1×10^4	7.34 $\times 10^{-13}$	0.342	0.0475	0.0554	0.0839
2×10^4	2.80 $\times 10^{-13}$	0.209	0.0198	0.0386	0.0943
4×10^4	1.22 $\times 10^{-13}$	0.194	0.0123	0.0341	0.0955
8×10^4	1.12 $\times 10^{-13}$	0.101	0.0066	0.0261	0.0924

Real Data Experiment - Microsoft Academic Graph:

- ▶ MAG1 ($N = 37680$, $K = 3$); MAG2 ($N = 19457$, $K = 3$); $L = 10$

Datasets	Proposed		GeoNMF		CD-MVSI	
	Avg. SRC	Time(s.)	Avg. SRC	Time(s.)	Avg. SRC	Time(s.)
MAG1	0.125	0.26	0.122	1.79	0.089	0.59
MAG2	0.441	0.23	0.240	4.66	0.249	0.53

References

- ▶ Xueyu Mao, Purnamrita Sarkar, and Deepayan Chakrabarti. On mixed memberships and symmetric nonnegative matrix factorizations. In International Conference on Machine Learning, pages 2324-2333, 2017.
- ▶ Kejun Huang and Xiao Fu. Detecting overlapping and correlated communities without pure nodes: Identifiability and algorithm. In International Conference on Machine Learning, pages 2859-2868, 2019.