

Detecting Overlapping and Correlated Communities without Pure Nodes: Identifiability and Algorithm

Kejun Huang

University of Florida

Xiao Fu

Oregon State University

International Conference on Machine Learning 2019

Mixed-membership Stochastic Blockmodel

MMSB [Airoldi et al., 2008]

- ▶ Given a graph adjacency matrix A
- ▶ An edge is present/absent follows Bernoulli

$$\Pr(A_{ij} = \{0, 1\}) = P_{ij}^{A_{ij}} (1 - P_{ij})^{1-A_{ij}}$$

- ▶ $P = M^T B M$: $B \in [0, 1]^{k \times k}$ community interaction
 $m_i \in \Delta = \{x : x \geq 0, \mathbf{1}^T x = 1\}$ mixed-membership of node i
- ★ Task: Uniquely identify (part of) M from data A
- ★ Challenges: **identifiability & scalability**

2nd-order Graph Moment

inspired by Anandkumar et al. [2014]

- ▶ Divide the network into three sets of nodes \mathcal{S}_0 , \mathcal{S}_1 , and \mathcal{S}_2
 - \mathcal{S}_2 : n nodes interested in finding their memberships
 - \mathcal{S}_1 : $k - 1$ nodes
 - \mathcal{S}_0 : all the other nodes to act as 2-star samples

$$\text{▶ } \hat{Y}_{i_1 i_2} = \frac{1}{|\mathcal{S}_0|} \sum_{i_0 \in \mathcal{S}_0} A_{i_0 i_1} A_{i_0 i_2} \quad i_1 \in \mathcal{S}_1 \quad i_2 \in \mathcal{S}_2$$

$$\text{▶ } Y_{i_1 i_2} = \mathbb{E}[\hat{Y}_{i_1 i_2}] = \mathbf{m}_{i_1}^\top \mathbf{B}^\top \left(\frac{1}{|\mathcal{S}_0|} \sum_{i_0 \in \mathcal{S}_0} \mathbf{m}_{i_0} \mathbf{m}_{i_0}^\top \right) \mathbf{B} \mathbf{m}_{i_2}$$

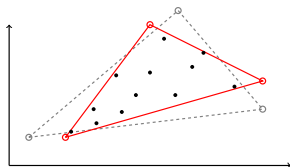
$$\text{▶ Let } \Sigma = \mathbb{E}[\mathbf{m}_{i_0} \mathbf{m}_{i_0}^\top] \text{ and } |\mathcal{S}_0| \rightarrow \infty, \text{ then } \hat{Y} \rightarrow \mathbf{M}_1^\top \mathbf{B}^\top \Sigma \mathbf{B} \mathbf{M}_2$$

$$\mathbf{Y} = \Xi \mathbf{M}_2$$

★ Can we uniquely recover $\mathbf{M}_2 \in \Delta^n$ from $\mathbf{Y} \in \mathbb{R}^{(k-1) \times n}$?

Geometric Interpretation

$$\mathbf{y}_{i_2} = \mathbf{\Xi} \mathbf{m}_{i_2} = \sum_{j=1}^k \xi_j m_{ji_2} \quad \mathbf{m}_{i_2} \in \Delta$$



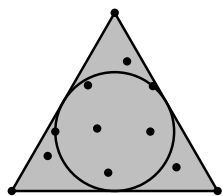
- ▶ \mathbf{y}_{i_2} is a **convex combination** of ξ_1, \dots, ξ_k
- ▶ \mathbf{y}_{i_2} belongs to the **convex hull** of ξ_1, \dots, ξ_k
- ▶ There are infinitely many enclosing simplexes
- ★ Intuition: Find the one with **minimum volume**

$$\begin{aligned} & \underset{\mathbf{\Xi}, \mathbf{M}_2}{\text{minimize}} \quad \frac{1}{(k-1)!} \left| \det \left[\begin{array}{ccc} \xi_1 - \xi_k & \cdots & \xi_{k-1} - \xi_k \end{array} \right] \right| \\ & \text{subject to} \quad \mathbf{Y} = \mathbf{\Xi} \mathbf{M}_2, \quad \mathbf{M}_2 \geq 0, \quad \mathbf{I}^\top \mathbf{M}_2 = \mathbf{1}. \end{aligned}$$

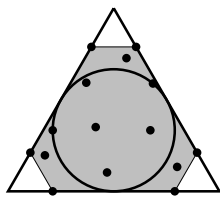
Definition: Sufficiently Scattered (informal)

Let \mathcal{D} be a “hyper-disc” on the hyperplane $\mathbf{I}^\top \mathbf{x} = 1$ defined as $\mathcal{D} = \{\mathbf{x} \in \mathbb{R}^k : \|\mathbf{x}\|^2 \leq \frac{1}{k-1}, \mathbf{I}^\top \mathbf{x} = 1\}$. A matrix \mathbf{M} , with all its columns in Δ , is called **sufficiently scattered** if $\mathcal{D} \subseteq \text{conv}(\mathbf{M})$.

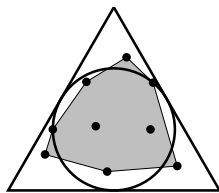
[Huang et al., 2014, 2016, 2018]



Pure node



Sufficiently scattered



Not identifiable

► Equivalently, define $\tilde{Y} = \begin{bmatrix} Y \\ I^\top \end{bmatrix}$, $\tilde{\Xi} = \begin{bmatrix} \Xi \\ I^\top \end{bmatrix}$,

$$\begin{aligned} & \underset{\tilde{\Xi}, M_2}{\text{minimize}} && \left| \det \tilde{\Xi} \right| \\ & \text{subject to} && \tilde{Y} = \tilde{\Xi} M_2, M_2 \geq 0, e_k^\top \tilde{\Xi} = I^\top. \end{aligned} \quad (\$)$$

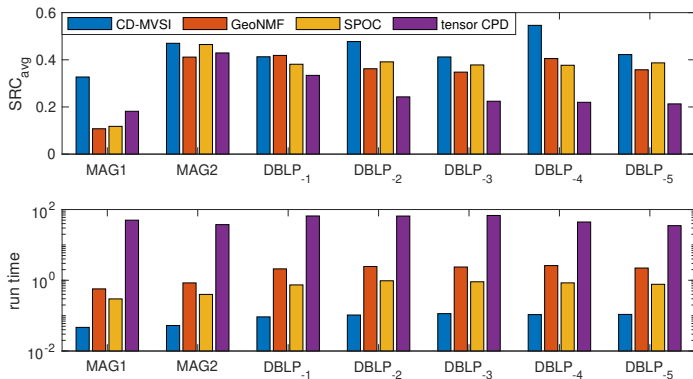
Theorem [Fu et al., 2015, Lin et al., 2015]

Suppose $Y = \Xi^{\natural} M_2^{\natural}$, where $\text{rank}(\tilde{\Xi}^{\natural}) = k$ and $M_2^{\natural} \in \Delta^n$ is **sufficiently scattered**. Let (M_*, Ξ_*) be an optimal solution for (\$\$), then there exists a permutation matrix $\Pi \in \mathbb{R}^{k \times k}$ such that

$$M_2^{\natural} = \Pi M_*, \quad \tilde{\Xi}^{\natural} = \Xi_* \Pi^\top.$$

Experiment

- ▶ Data sets:
 - Coauthorship data from Microsoft Academic Graph (MAG) and DBLP [Mao et al., 2017]
 - Groundtruth community: “field of study” in MAG and venues in DBLP



References I

- Edoardo M Airoldi, David M Blei, Stephen E Fienberg, and Eric P Xing. Mixed membership stochastic blockmodels. *Journal of Machine Learning Research*, 9: 1981–2014, 2008.
- Animashree Anandkumar, Rong Ge, Daniel Hsu, and Sham M Kakade. A tensor approach to learning mixed membership community models. *Journal of Machine Learning Research*, 15(1):2239–2312, 2014.
- Xiao Fu, Wing-Kin Ma, Kejun Huang, and Nicholas D Sidiropoulos. Blind separation of quasi-stationary sources: Exploiting convex geometry in covariance domain. *IEEE Transactions on Signal Processing*, 63(9), 2015.
- Kejun Huang, Nicholas D Sidiropoulos, and Ananthram Swami. Non-negative matrix factorization revisited: Uniqueness and algorithm for symmetric decomposition. *IEEE Transactions on Signal Processing*, 62(1):211–224, 2014.
- Kejun Huang, Xiao Fu, and Nikolaos D Sidiropoulos. Anchor-free correlated topic modeling: Identifiability and algorithm. In *Advances in Neural Information Processing Systems*, pages 1786–1794, 2016.
- Kejun Huang, Xiao Fu, and Nicholas Sidiropoulos. Learning hidden Markov models from pairwise co-occurrences with application to topic modeling. In *International Conference on Machine Learning*, pages 2068–2077. PMLR, 2018.

References II

- Chia-Hsiang Lin, Wing-Kin Ma, Wei-Chiang Li, Chong-Yung Chi, and ArulMurugan Ambikapathi. Identifiability of the simplex volume minimization criterion for blind hyperspectral unmixing: The no-pure-pixel case. *IEEE Transactions on Geoscience and Remote Sensing*, 53(10):5530–5546, 2015.
- 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.
- Krzysztof Nowicki and Tom A B Snijders. Estimation and prediction for stochastic blockstructures. *Journal of the American Statistical Association*, 96(455): 1077–1087, 2001.
- Tom AB Snijders and Krzysztof Nowicki. Estimation and prediction for stochastic blockmodels for graphs with latent block structure. *Journal of Classification*, 14(1): 75–100, 1997.