

Network Community Detection Using Higher Order Interactions

Xianshi Yu

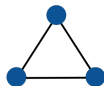
Joint work with Ji Zhu

Aug 03, 2020

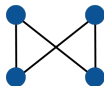
Table of contents

- 1 Real Data Observation
- 2 Ideas
- 3 Method
- 4 Experiment on lawyers' co-work network
- 5 Model & Theory

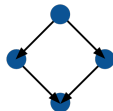
Real Data Observation



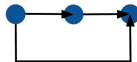
triangle



by-fan



by-parallel



feed-forward loop

Certain subgraphs are abundant. [4,5,6]

triangle: social network; world wide webs

by-fan: transcriptional gene regulation network; neural network

by-parallel: neural network; food web

...

Real Data Observation

The pattern of observed subgraphs could reflect the underlying community memberships.

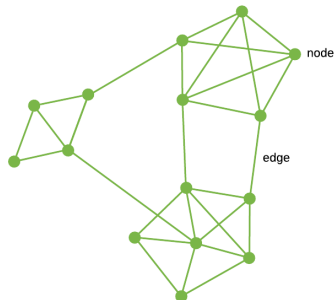
- How to exploit the configuration of Δ 's in community detection?
- How to model the abundance and distribution of Δ ?

Notations

- $\{1, 2, \dots, n\}$ index the nodes.
- $A_{n \times n}$ (adjacency matrix) describes edges of an observed network.
- $A_{ij} = 1$ if nodes i and j have an edge, $A_{ij} = 0$ otherwise.
- $\Delta_{n \times n \times n}$ represents triangles: $\Delta_{ijk} = 1$ iff nodes i, j and k share a triangle. $\Delta_{ijk} = 0$ otherwise.
- ★ $\Delta_{ijk} = A_{ij}A_{jk}A_{ki}$

Goal

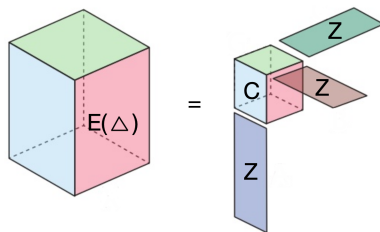
- $g_i \in \{1, 2, \dots, K\}$ is the community index of node i .
- $Z_{n \times K}$ is a matrix of all zeros, except that $Z_{ig_i} = 1$ for all i .



Consider the formation of triangles as depending only on community membership

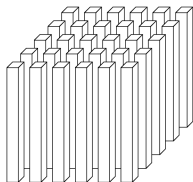
$$P(\Delta_{ijk} = 1) = C_{g_i, g_j, g_k}$$

i.e. $\mathbb{E}(\Delta_{n \times n \times n}) = C_{K \times K \times K} \times_1 Z \times_2 Z \times_3 Z$

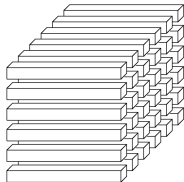


Definition of n -mode product ▶ ref

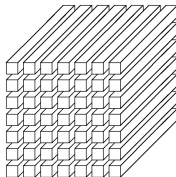
To obtain $C_{K \times K \times K} \times_1 Z_{n \times K}$
 Consider mode-1 fibers of C



(a) Mode-1 (column) fibers:
 $C_{\cdot jk}$



(b) Mode-2 (row) fibers:
 $C_{i \cdot k}$



(c) Mode-3 (tube) fibers:
 $C_{ij \cdot}$

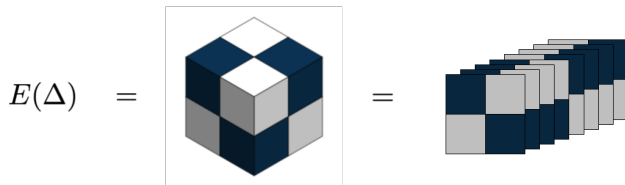
Calculate the products of each mode-1 fiber $C_{\cdot jk}$ and Z , namely $ZC_{\cdot jk}^T$.
 Arrange the resulted vectors accordingly to form an $n \times K \times K$ tensor.

Remarks

$$\mathbb{E}(\Delta) = C \times_1 Z \times_2 Z \times_3 Z \quad (1)$$

- In Stochastic Block Model (SBM), community structure is on $E(A)$.
- (1) does not define a generative model, but a 'constraint' on network models. SBM satisfies (1).

How to infer $g_i, i \in \{1, \dots, n\}$ from Δ ?

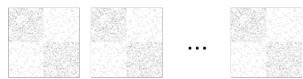




$$E(\Delta_{..k}), k \in G_1$$



$$E(\Delta_{..k}), k \in G_2$$



$$\sum_{k \in G_1} \Delta_{..k}$$



$$\sum_{k \in G_2} \Delta_{..k}$$

- $\mathbb{E}(\Delta)$ is a block tensor, and slices of it, i.e $\mathbb{E}(\Delta_{..k})$, are block matrices. Spectral decomposition of $\mathbb{E}(\Delta_{..k})$ reveals g_j .
- $\Delta_{..k}$ is likely very sparse.

Method

Algorithm

Input: $A_{n \times n}, K$

Output: \hat{Z}

Step 1. Calculate $\Delta = (A_{i,j}A_{j,k}A_{i,k})_{i,j,k}$

Step 2. Obtain an initial estimate $Z_{n \times K}^0$

Step 3. Calculate **sums of slices** of $\Delta_{n \times n \times n}$ w.r.t groups defined by Z^0 .
i.e. $S_{\Delta}^l := \Delta \times_3 Z_{\cdot,l}^0$, for $l \in \{1, \dots, K\}$.

Step 4. Apply **synchronized spectral decomposition** to S_{Δ}^l ,
 $l \in \{1, \dots, K\}$ to obtain $\hat{U}_{n \times K}$

Step 5. Perform K-means on \hat{U} to obtain \hat{Z}

Method—synchronized spectral decomposition (SynSD)

$$\hat{U} = \arg \max_{U_{n \times K}^T U = I} \sum_{l=1}^K \|U^T S_{\Delta}^l U\|_F^2$$

- SynSD finds ‘shared singular vectors’.
- SynSD is a Grassmann manifold optimization problem, for which ample literature and packages are available.

Experiment on lawyers' co-work network

71 attorneys (1104 edges)

seniority	status	gender	office	age
range:[1,32] median:7	partner:36 associate:35	man:53 woman:18	Boston:48 Hartford:19 Providence:4	range:[26,67] median:39

practice	law school
litigation:41 corporate:30	harvard, yale:15 ucon:28 other:28

Experiment on lawyers' co-work network

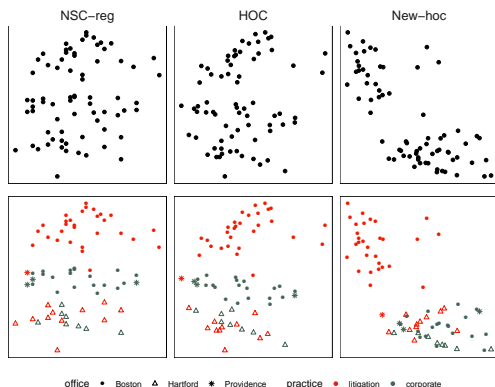
Benchmarks

- Normalized Spectral Clustering with regularization (NSC-reg)
- High-Order Clustering [1, 2] (HOC)
 - perform NSC on $\sum_{k=1}^n \Delta_{..k}$

Initial estimate Z^0 of the new method

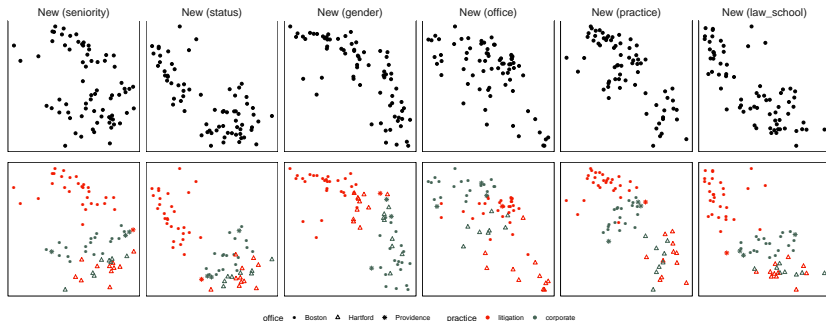
- HOC
- Attribute of lawyers

Experiment on lawyers' co-work network



- Column 1&2: leading two eigenvectors of NSC-reg and HOC
- Column 3: \hat{U} of the new method ($K = 2$)

Experiment on lawyers' co-work network



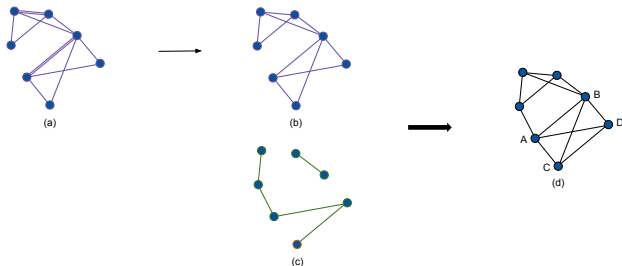
\hat{U} obtained when attributes in bracket are taken as Z^0

Model & Theory

A model that satisfies $\mathbb{E}(\Delta) = C \times_1 Z \times_2 Z \times_3 Z$

— triangle mechanism + edge mechanism

- $T \sim$ independent Bernoulli with $\ddot{P}_{n \times n \times n} = \ddot{B} \times_1 Z \times_2 Z \times_3 Z$
- $\ddot{A} \sim$ independent Bernoulli with $\ddot{P}_{n \times n}$
- $A_{i,j} = \max\{\ddot{A}_{i,j}, \mathbf{1}(\sum_k T_{i,j,k} > 0)\}$



Conditions

- $C_1 \ln^4 n \leq \ddot{d} \leq C_2 n^{\frac{2}{5}}$, where $\ddot{d} = n^2 \ddot{P}_{\max}$
- $\ddot{d} \leq C_3 n^{\frac{1}{3}} \ddot{d}^{\frac{1}{6}}$, where $\ddot{d} = n \ddot{P}_{\max}$

Notations

- $\sigma_{\min}^l = \sigma_K(\mathbb{E}(S_{\Delta}^l))$, $\sigma_{\max}^l = \sigma_1(\mathbb{E}(S_{\Delta}^l))$

Theorem 1

With $n > N$ and S_{Δ}^l calculated from a fixed Z^0 , if the conditions above are satisfied, and \hat{U} is the global optimum of synchronized spectral decomposition, then with probability at least $1 - n^{-r}$,

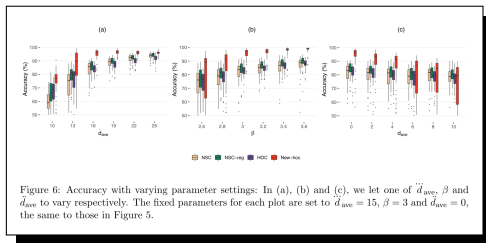
$$\frac{|\mathcal{M}|}{n} \leq C_4 K \frac{K \ddot{d} + \ddot{d}^{\frac{1}{2}} \sum_{l=1}^K \sigma_{\max}^l}{\sum_{l=1}^K (\sigma_{\min}^l)^2}.$$

Here \mathcal{M} is the set of misclustered nodes and N, r, C_1, \dots, C_4 are absolute constants.

Remark

In the simple case of the model where $\ddot{P} = 0$, consider n growing, if Z^0 is independent of A with a confusion matrix proportionally constant and \ddot{B}/\ddot{P}_{\max} constant, then the upper bound in Theorem 1 is $O(\ddot{d}^{-1/2})$.

Simulation results are in the manuscript



Conclusion

- A new method for network community detection that takes as input the observed \triangle 's.
- The new method tries to explain edges in networks not only by affiliation but also by roles in higher order interaction.
- Consistency theory that involves analyzing dependent objects.

Future work:

Other subgraphs, e.g. \bowtie , \diamond , \vee .

Previous methods that use triangles

- 1 Benson, A. R., Gleich, D. F. and Leskovec, J. (2016). Higher-order organization of complex networks. *Science*, 353(6295), 163-166.
- 2 Paul, S., Milenkovic, O. and Chen, Y. (2018). Higher-Order Spectral Clustering under Superimposed Stochastic Block Model. *arXiv preprint arXiv:1812.06515*.
- 3 Vandecappelle, M., Boussé, M., Van Eeghem, F. and De Lathauwer, L. (2016). Tensor decompositions for graph clustering. Internal Report, (16-170).

Abundance of subgraphs in real world networks.

- 4 Milo, R., Shen-Orr, S., Itzkovitz, S., Kashtan, N., Chklovskii, D. and Alon, U. (2002). Network motifs: simple building blocks of complex networks. *Science*, 298(5594), 824-827.
- 5 Mangan, S. and Alon, U. (2003). Structure and function of the feed-forward loop network motif. *Proceedings of the National Academy of Sciences*, 100(21), 11980-11985.
- 6 Yaveroglu, Ö. N., Malod-Dognin, N., Davis, D., Levnajic, Z., Janjic, V., Karapandza, R., ... and Pržulj, N. (2014). Revealing the hidden language of complex networks. *Scientific reports*, 4, 4547.