



Pairwise covariates-adjusted block model for community detection

Yang Feng Department of Statistics, Columbia University Many systems can be represented as networks.



5

- online social network, friendship between users
- citation networks
- Hyperlinks between web pages
- email networks

Note: images credit to crmswitch.com and www.kai-arzheimer.com/networkpics/mutualworld.png.

Community detection (Fortunato, 2010; Malliaros et al, 2013): Identifying organization of vertices in clusters, based on the network topology.



community structures detected (Girvan & Newman, 2003) in a web network

Communities are functional or structural units in a networked system

- social groupings in a social network
- related papers on a single topic in a citation network
- customers of similar interest in a network of purchase relationships

...

1

- A graph G(V, E) with vertex set V and edge set E
- ▶ The adjacency matrix $A \in \{0,1\}^{n \times n}$, where $A_{ij} = 1$ if and only if $(i,j) \in E$.

Stochastic Block Model (Holland, 1983)

n nodes, K non-overlapping communities, membership c = {c_i}ⁿ_{i=1}, connectivity matrix B ∈ [0, 1]^{K×K}

1

- $P(A_{ij} = 1) = B_{c_i c_j}$ $(1 \le i < j \le n)$
- Extensions: Degree-Corrected Stochastic Block Model (Karrer & Newman, 2010), Mixed Membership Stochastic Block Model (Airoldi et al., 2009).



Networks often appear with additional covariate information

- online social networks with users' personal profile information
- citation networks with papers' authors, keywords and abstracts

Motivations

. . .

- 1. more accurate community detection
- 2. identify relevant covariates
- 3. interplay between node and edge information

Some existing methods: Nallapati and Cohen, 2008; Chang and Blei, 2010; Akoglu et al., 2012; Ruan et al., 2013; Yang et al., 2013; Zhang et al, 2013; Binkiewicz et al, 2015; etc. Mostly heuristic or based on algorithms, not much model based work. Networks often appear with additional covariate information

- online social networks with users' personal profile information
- citation networks with papers' authors, keywords and abstracts

▶ ...

Motivations

- 1. more accurate community detection
- 2. identify relevant covariates
- 3. interplay between node and edge information

Some existing methods: Nallapati and Cohen, 2008; Chang and Blei, 2010; Akoglu et al., 2012; Ruan et al., 2013; Yang et al., 2013; Zhang et al, 2013; Binkiewicz et al, 2015; etc. Mostly heuristic or based on algorithms, not much model based work.

Networks often appear with additional covariate information

- online social networks with users' personal profile information
- citation networks with papers' authors, keywords and abstracts

▶ ...

Motivations

- 1. more accurate community detection
- 2. identify relevant covariates
- 3. interplay between node and edge information

Some existing methods: Nallapati and Cohen, 2008; Chang and Blei, 2010; Akoglu et al., 2012; Ruan et al., 2013; Yang et al., 2013; Zhang et al, 2013; Binkiewicz et al, 2015; etc. Mostly heuristic or based on algorithms, not much model based work.



Suppose, in addition to the adjacency matrix A, the available covariates information is formulated in a covariate matrix $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T \in \mathbb{R}^{n \times d}$, where $\mathbf{x}_i \in \mathbb{R}^d$ is the *i*-th node's covariate vector.

Goal: estimate \mathbf{c} from the observations A and X.

Question: How to integrate the nodal information with the adjacency matrix?



Suppose, in addition to the adjacency matrix A, the available covariates information is formulated in a covariate matrix $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T \in \mathbb{R}^{n \times d}$, where $\mathbf{x}_i \in \mathbb{R}^d$ is the *i*-th node's covariate vector.

Goal: estimate \mathbf{c} from the observations A and X.

Question: How to integrate the nodal information with the adjacency matrix?



Suppose, in addition to the adjacency matrix A, the available covariates information is formulated in a covariate matrix $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T \in \mathbb{R}^{n \times d}$, where $\mathbf{x}_i \in \mathbb{R}^d$ is the *i*-th node's covariate vector.

Goal: estimate \mathbf{c} from the observations A and X.

Question: How to integrate the nodal information with the adjacency matrix?

Two possible model assumptions

- Assuming both A and X depend on c.
- Example: citation network; c: research topic, X: keywords.



- ► Assuming A depends on both **c** and X.
- Example: high-school friendship network; c: grade, X: ethnicity.



Two possible model assumptions

- Assuming both A and X depend on c.
- Example: citation network; c: research topic, X: keywords.



- ► Assuming A depends on both **c** and X.
- Example: high-school friendship network; c: grade, X: ethnicity.



First model: a generic two-step inference procedure (Weng and F., 2017)

Assuming $A \perp X \mid \mathbf{c}$,

(1) Parameter estimation based on $P(A, X; \theta)$:

$$P(A, X; \theta) = \sum_{\mathbf{c}} P(A \mid \mathbf{c}) P(X \mid \mathbf{c}) P(\mathbf{c})$$
(1)
= $P(X; \theta_1) \sum_{\mathbf{c}} P(A \mid \mathbf{c}; \theta_2) P(\mathbf{c} \mid X; \theta_3),$ (2)

50

where $\theta = (\theta_1, \theta_2, \theta_3)$ indexes a family of generative models (here, we are not assuming specific parametric forms)

(2) Posterior inference according to $P(\mathbf{c} \mid A, X)$.

$$P(\mathbf{c} \mid A, X) = \frac{P(A \mid \mathbf{c})P(X \mid \mathbf{c})P(\mathbf{c})}{\sum_{\mathbf{c}} P(A \mid \mathbf{c})P(X \mid \mathbf{c})P(\mathbf{c})} = \frac{P(A \mid \mathbf{c})P(\mathbf{c} \mid X)P(X)}{\sum_{\mathbf{c}} P(A \mid \mathbf{c})P(\mathbf{c} \mid X)}$$
$$= \frac{P(A \mid \mathbf{c})P(\mathbf{c} \mid X)}{\sum_{\mathbf{c}} P(A \mid \mathbf{c})P(\mathbf{c} \mid X)}.$$
(3)



(a) SBM:

$$\mathcal{P}(A \mid \mathbf{c}) = \prod_{i < j} B^{A_{ij}}_{c_i c_j} (1 - B_{c_i c_j})^{1 - A_{ij}}$$

(b) multi-logistic regression:

$$P(\mathbf{c} \mid X) = \prod_{i} \frac{\exp(\beta_{c_{i}}^{T} \mathbf{x}_{i})}{\sum_{k=1}^{K} \exp(\beta_{k}^{T} \mathbf{x}_{i})},$$

where $B = (B_{ab}) \in [0, 1]^{K \times K}$, $\beta = (\beta_1^T, \dots, \beta_K^T)^T \in \mathbb{R}^{K_p}$. Here, $\beta_K = 0$ for identifiability.

Community detection consistency

$$\begin{array}{ll} (\textit{strong consistency}) & P(\hat{\mathbf{c}} = \mathbf{c}) \to 1, \text{ as } n \to \infty, \\ (\textit{weak consistency}) & P\Big(\frac{1}{n}\sum_{i=1}^n \mathbf{1}(\hat{c}_i \neq c_i) < \epsilon\Big) \to 1, \text{ as } n \to \infty, \forall \epsilon > 0. \end{array}$$

Here, $B = \rho_n \overline{B}$ with \overline{B} fixed and sparsity level $\rho_n = P(A_{ij} = 1) \rightarrow 0$ as $n \rightarrow \infty$ (Bickel and Chen (2009); Zhao et al. (2012); Bickel et al. (2013))

Note: both versions of consistency are interpreted upon a permutation of the community labels.

Marginalization over **c**.

$$(\hat{\beta}, \hat{B}) = \arg \max_{\substack{\beta_{K}=0, \ \beta \in \mathbb{R}^{K_{p}} \\ B \in [0,1]^{K \times K}, B^{T}=B}} \sum_{\mathbf{c}} \prod_{i < j} B_{c_{i}c_{j}}^{A_{ij}} (1 - B_{c_{i}c_{j}})^{1 - A_{ij}} \cdot \prod_{i} \frac{e^{\beta_{c_{i}}^{-} \mathbf{x}_{i}}}{\sum_{k=1}^{K} e^{\beta_{k}^{T} \mathbf{x}_{i}}}, (4)$$

$$\hat{\mathbf{c}} = \arg \max_{\mathbf{c} \in \{1, ..., K\}^{n}} \prod_{i < j} \hat{B}_{c_{i}c_{j}}^{A_{ij}} (1 - \hat{B}_{c_{i}c_{j}})^{1 - A_{ij}} \cdot \prod_{i} \frac{e^{\hat{\beta}_{c_{i}}^{-T} \mathbf{x}_{i}}}{\sum_{k=1}^{K} e^{\hat{\beta}_{k}^{T} \mathbf{x}_{i}}}.$$

$$(5)$$

Ŧ

Condition 1

 \bar{B} has no two identical columns.

Condition 2

 $(c_1, \mathbf{x}_1), \ldots, (c_n, \mathbf{x}_n) \stackrel{iid}{\sim} (c, \mathbf{x})$ with $\mathbb{E}(\mathbf{x}\mathbf{x}^T) \succ 0$, where $\succ 0$ represents the matrix being positive definite.

0

Condition 3

There exist constants κ_1 and κ_2 such that for sufficiently large *t*, we have

$$P(\|\boldsymbol{x}\|_2 > t) \leq \kappa_1 e^{-\kappa_2 t}.$$

Theorem 1

Assume the data (A, X) follows NSBM and Conditions 1-3 hold. In addition, assume $\frac{n\rho_n}{\log n} \to \infty$ as $n \to \infty$. Then, we have as $n \to \infty$

$$P(\hat{\boldsymbol{c}} = \boldsymbol{c})
ightarrow 1, \quad \sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \stackrel{d}{
ightarrow} N(\boldsymbol{0}, I^{-1}(\boldsymbol{\beta})),$$

where $I(\beta)$ is the Fisher information for the multi-logistic regression problem of regressing c on X.

Remark: condition on expected degree the same as Bickel and Chen (2009) and Zhao et al. (2012).

Remark: Computational infeasible even to evaluate the objective function at any non-degenerate point.

Theorem 1

Assume the data (A, X) follows NSBM and Conditions 1-3 hold. In addition, assume $\frac{n\rho_n}{\log n} \to \infty$ as $n \to \infty$. Then, we have as $n \to \infty$

$$P(\hat{\boldsymbol{c}} = \boldsymbol{c})
ightarrow 1, \quad \sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \stackrel{d}{
ightarrow} N(\boldsymbol{0}, I^{-1}(\boldsymbol{\beta})),$$

where $I(\beta)$ is the Fisher information for the multi-logistic regression problem of regressing c on X.

Remark: condition on expected degree the same as Bickel and Chen (2009) and Zhao et al. (2012).

Remark: Computational infeasible even to evaluate the objective function at any non-degenerate point.



More in: Weng and **F.** (2016), Community detection with nodal information, https://arxiv.org/abs/1610.09735.

- Variational EM
- Profile likelihood
- Semi-definite Programming

Assuming A depends on both \mathbf{c} and X.



(h)

- ▶ Given X_{n×p}, define a p-dimensional covariate Z_{ij} = g(X_i, X_j). It is a measure of similarity of nodes i and j in different aspects.
- ► For example, |X_{i1} X_{j1}| measures the distance of the covariates for two nodes.

50

▶ Given community assignment c and pairwise covariates Z, assume A_{ij} are independent Poisson distributed with rate

$$\lambda_{ij} := P(A_{ij} = 1) = B_{c_i c_j} e^{Z_{ij}^T \gamma^0}$$

- Adjustment factor: $\exp\{Z_{ij}^T \gamma^0\}$.
- Reduces to the standard SBM when $\gamma^0 = 0$.
- Let $B = \rho_n \overline{B}$ with \overline{B} fixed and $\rho_n \to 0$.

► The likelihood function of the complete graph model (CGM) is defined as

$$f(A, e|Z, \gamma, B, \pi) = \prod_{i=1}^{n} \pi_{e_i} \prod_{i < j} B_{e_i e_j}^{A_{ij}} e^{A_{ij} Z_{ij}^T \gamma} e^{-B_{e_i e_j} e^{Z_{ij}^T \gamma}},$$

where $\pi_a = P(e = a)$.

We then define the following maximum likelihood estimators for graph model (GM):

$$\hat{\gamma} = \arg \max_{\gamma \in \Gamma} \sum_{e} f(A, e | Z, \gamma), \ \hat{c} = \arg \max_{e} f(A, e | Z, \hat{\gamma}),$$

where Γ is a compact set in \mathbb{R}^{p} . Basically, GM is the community marginalization of CGM.



$$n_k(\mathbf{e}) = \sum_{i=1}^n \mathbf{1}_{e_i=k}$$
, and $s_{\mathbf{e}}(k, l) = \{(i, j) | e_i = k, e_j = l, i \neq j\}$.

Condition 1 (non-degeneracy)

 $n_k(\mathbf{e}) \to \infty$ for any k.

Condition 2 (non-collinearity)

There exists a community pair (k, l) s.t. $|s_e(k, l)| \ge p + 1$ and the matrix $[(\mathbf{Z}_{ij}^T \mathbf{1})]_{(i,j)\in s_e(k,l)} \in \mathbb{R}^{|s_e(k,l)| \times (p+1)}$ is full column rank.

For any given e, let

$$ilde{\gamma}(e) = rg\max_{\gamma\in \Gamma} f(A, e|Z, \gamma).$$

Theorem 3

Under Conditions 1 and 2, as $n \to \infty$, $\tilde{\gamma}(e) \xrightarrow{p} \gamma^0$.

Now, let

$$\hat{c} = rg\max_{e'} f(A, e' | Z, ilde{\gamma}(e)).$$

Theorem 4

Under Conditions 1 and 2, as $n \to \infty$, if $N\rho_n \to \infty$, the asymptotical distribution of $\sqrt{N\rho_n}(\tilde{\gamma}(\mathbf{e}) - \gamma^0)$ is multivariate normal distribution with mean **0** and covariance matrix $I^{-1}(\gamma^0)$, where $I(\gamma^0) = \sum_{ab} \bar{B}_{ab} \pi_a \pi_b (\nu_2 - \nu_0^{-1} \nu_1 \nu_1^T)$.

For any given e, let

$$ilde{\gamma}(e) = rg\max_{\gamma\in \Gamma} f(A, e|Z, \gamma).$$

Theorem 3

Under Conditions 1 and 2, as $n \to \infty$, $\tilde{\gamma}(e) \xrightarrow{p} \gamma^0$.

Now, let

$$\hat{c} = \arg \max_{e'} f(A, e' | Z, \tilde{\gamma}(e)).$$

Theorem 4

Under Conditions 1 and 2, as $n \to \infty$, if $N\rho_n \to \infty$, the asymptotical distribution of $\sqrt{N\rho_n}(\tilde{\gamma}(\mathbf{e}) - \gamma^0)$ is multivariate normal distribution with mean **0** and covariance matrix $I^{-1}(\gamma^0)$, where $I(\gamma^0) = \sum_{ab} \bar{B}_{ab} \pi_a \pi_b (\nu_2 - \nu_0^{-1} \nu_1 \nu_1^T)$.

Theorem 5

Under PCABM, when $\tilde{\gamma}(\mathbf{e})$ is consistent, $\tilde{\mathbf{c}}(\mathbf{e})$ is weakly consistent if $\varphi_n \to \infty$ and strongly consistent if $\varphi_n / \log n \to \infty$.

Spectral clustering with adjustment (SCWA)

Observation

$$E\left[\frac{A_{ij}}{\exp\{Z_{ij}^{T}\gamma^{0}\}}\right] = B_{c_{i}c_{j}}.$$

Intuitively, we would like to perform spectral clustering on $A' = \{A'_{ij}\}$, where

$$\mathsf{A}'_{ij} = \frac{\mathsf{A}_{ij}}{\exp\{Z_{ij}^T \tilde{\gamma}\}}.$$

Remark

We need to choose a lower bound v(n) for the adjusted factor, to control large variances of the adjusted variable. When $e^{\mathbf{Z}_{ij}^T \tilde{\gamma}} < v(n)$, we use v(n) as the adjustment factor.

Theorem 6 (Spectral bound of poisson random matrices).

Let A be the adjacency matrix generated by PCABM (M, B, Z, γ^0) , and the adjusted adjacency matrix A' is derived with some proper threshold value v(n) > 0. Assume that $nP_{\max} \le d$ and $d \ge C_1 \log n$ for $C_1 > 0$. Then, for any r > 0, there exists a constant C such that

$$|A'-P\| \leq C\sqrt{d}/v(n)$$

with probability at least $1 - n^{-r \times v(n)}$.

We define two measure of estimation error, an overall relative error and a worst case relative error:

$$L_{1}(\hat{M}, M) = n^{-1} \min_{Q \in S_{\kappa}} \|\hat{M}Q - M\|_{0},$$

$$L_{2}(\hat{M}, M) = \min_{Q \in S_{\kappa}} \max_{1 \le k \le \kappa} n_{k}^{-1} \|(\hat{M}Q)_{G_{k}} - M_{G_{k}} \|_{0},$$

Theory (Cont')

Theorem 7 (Error bound)

Let A be an adjacency matrix generated from a PCABM (M, B, Z, γ^0) , where $B = \rho_n \overline{B}$ for some $\rho_n \ge C_1 \log n/n$ and with $\overline{B}'s$ minimum absolute eigenvalue bounded below by $\tau > 0$ and $\max_{kl} \overline{B}(k, l) = 1$. The adjusted adjacency matrix A' is derived with some threshold value v(n) > 0. Let \hat{M} be the output of spectral clustering using $(1 + \epsilon)$ approximate k-means. For any r > 0, there exists an absolute constant C such that if

$$(2+\epsilon)\frac{Kn}{n_{\min}^2\tau^2\rho_nv^2(n)} < C,$$

then with probability at least $1 - n^{-r \times v(n)}$,

$$\mathcal{L}_2(\hat{M},M) \leq C^{-1}(2+\epsilon) rac{Kn}{n_{\min}^2 au^2
ho_n v^2(n)},$$

$$L_1(\hat{M}, M) \leq C^{-1}(2+\epsilon) \frac{Kn'_{\max}}{n_{\min}^2 \tau^2 \rho_n v^2(n)}$$

Pairwise covariates-adjusted block model for community detection, IMS workshop at NUS 4 ロト 4 昼 ト 4 匡 ト 4 匡 ト 3 三 の 9 0 25



- $\blacktriangleright B_0 = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}.$
- $\blacktriangleright \ n = 100, 200, 300, 400, 500.$
- ▶ $\rho_n = 2(\log n)^{1.5}/n.$
- $\gamma^0 = (0.4, 0.8, 1.2, 1.6, 2).$
- ► The components of X_i are i.i.d. and follow the distributions: B(1,0.1), Pois(0.1), U(0,1), Exp(0.3), N(0,0.3).
- The pairwise covariates $\mathbf{Z}_{ij} = |\mathbf{X}_i \mathbf{X}_j|$

We use two quantitative measures for evaluating the community detection performance.

Normalized Mutual Information (Ana and Jain, 2003):

$$\mathsf{NMI} = \frac{-2\sum_{i}\sum_{j}n_{ij}\log\left(\frac{n_{ij}\cdot n}{n_{i}\cdot n_{j}}\right)}{\sum_{i}n_{i}\cdot\log\left(\frac{n_{i}}{n}\right) + \sum_{j}n_{\cdot j}\log\left(\frac{n_{\cdot j}}{n}\right)}.$$

Adjusted Rand Index (Hubert and Arabie, 1985):

$$\mathsf{ARI} = \frac{\sum_{ij} \binom{n_{ij}}{2} - \frac{\sum_{i} \binom{n_{i}}{2} \sum_{j} \binom{n_{j}}{2}}{\binom{n}{2}}}{\frac{1}{2} \sum_{i} \binom{n_{i}}{2} + \frac{1}{2} \sum_{j} \binom{n_{i}}{2} - \frac{\sum_{i} \binom{n_{i}}{2} \sum_{j} \binom{n_{ij}}{2}}{\binom{n}{2}}},$$

where n_i . denotes the true number of nodes in community *i*, $n_{\cdot j}$ represents the number of nodes in the estimated community *j* and n_{ij} is the number of nodes belonging to community *i* but estimated to be in community *j*. In [0,1].



- SBM.MLE: MLE for SBM
- SBM.SC: Spectral clustering for SBM
- ▶ PCABM.MLE0: MLE for PCABM with random initialization
- PCABM.SCWA: Spectral clustering with adjustment for PCABM
- PCABMw.MLE0: MLE for PCABM with random initialization and misspecified pairwise function

ab)



DCBM setting



Simulation results under DCBM for different parameter settings

- The nodes are blogs about US politics and edges represent hyperlinks between them.
- 1,222 nodes and 16,714 edges.
- We let $Z_{ij} = \log(d_i \times d_j)$, where d_i is the degree for the *i*-th node.

	ARI	NMI	Errors
Karrer(2011)	-	0.72	-
Zhao(2012)	0.819	-	-
SCORE	0.819	0.725	58
PCABM.MLE	0.825	0.737	56

Performance comparison on political blogs data



- ▶ Network with 777 nodes and 4124 edges.
- Covariates: grade, gender, ethnicity and number of friends nominated.
- In each experiment, out of the three nodal covariates, school, ethnicity and gender, we viewed one covariate as the indicator for the underlying community, and performed community detection by using the pairwise covariates constructed with the other two covariates.
- ► For gender, school and ethnicity, we created dummy variables to represent whether the corresponding covariate values were the same for the pair of nodes. In addition, for number of nominated friends, we use Z_{ij} = log(n_i + 1) + log(n_j + 1).

Results



d

Ethnicity • White • Black

Pairwise covariates-adjusted block model for community detection, IMS workshop at NUS \triangleleft \square \blacktriangleright \triangleleft \vdash \triangleleft \vdash \triangleleft \vdash \triangleleft \equiv \flat \triangleleft \equiv \flat \parallel \equiv ϑ \triangleleft \triangleleft 23

	School	Race	Gender
SBM.SC	0.043	-0.024	0
SBM.MLE	0.048	0.138	-0.001
PCABM.MLE	0.894	0.914	0.03

- Combine the two models by dividing the covariates into two categories automatically.
- In the setting with high-dimensional covariates, penalized likelihood methods are more appealing for both community detection and variable selection.
- ► For very sparse networks, considering npn = O(1) seems to be a more realistic asymptotic framework. Under such asymptotics, community detection consistency is impossible.
- In this work, we assume the number of communities K is known. How to select K is an important problem in community detection. Some recent efforts towards this direction include Saldana, Yu and F. (2015); Le and Levina (2015); Wang and Bickel (2015); Lei (2016); Chen and Lei (2016); Li and Zhu (2017), etc.

Thank you very much!

References:

- Weng and F. (2016), Community detection with nodal information. https://arxiv.org/abs/1610.09735.
- Huang and F. (2018), Pairwise Covariates-adjusted Block Model for Community Detection, manuscript.

► SBM: (Daudin et al., 2008; Celisse et al., 2012; Bickel et al., 2013).

$$\max_{\boldsymbol{\theta}} \log P(A, X; \boldsymbol{\theta}) = \max_{\boldsymbol{\theta}, Q(\cdot)} \mathbb{E}_Q[\log P(A, X, \mathbf{c}; \boldsymbol{\theta}) - \log Q(\mathbf{c})], \tag{6}$$

where $Q(\cdot)$ denotes any joint distribution on **c**.

▶ Here, we consider the mean-field variational approach (Jordan et al., 1999),

$$\max_{\theta, Q \in \mathcal{Q}} \mathbb{E}_{Q}[\log P(A, X, \mathbf{c}; \theta) - \log Q(\mathbf{c})],$$
(7)

where $Q = \{Q : Q(\mathbf{c}) = \prod_{i=1}^{n} q_{ic_i}, \sum_k q_{ik} = 1, 1 \le i \le n\}$. The subset Q contains all the distributions under which \mathbf{c} is jointly independent.

• Denote the maximizer in (7) by $(\check{B},\check{\beta})$ and \check{c} .

SBM: (Daudin et al., 2008; Celisse et al., 2012; Bickel et al., 2013).

$$\max_{\boldsymbol{\theta}} \log P(\boldsymbol{A}, \boldsymbol{X}; \boldsymbol{\theta}) = \max_{\boldsymbol{\theta}, \boldsymbol{Q}(\cdot)} \mathbb{E}_{Q}[\log P(\boldsymbol{A}, \boldsymbol{X}, \mathbf{c}; \boldsymbol{\theta}) - \log Q(\mathbf{c})], \quad (6)$$

where $Q(\cdot)$ denotes any joint distribution on **c**.

► Here, we consider the mean-field variational approach (Jordan et al., 1999),

$$\max_{\boldsymbol{\theta}, \boldsymbol{Q} \in \boldsymbol{\mathcal{Q}}} \mathbb{E}_{Q}[\log P(A, X, \mathbf{c}; \boldsymbol{\theta}) - \log Q(\mathbf{c})],$$
(7)

where $Q = \{Q : Q(\mathbf{c}) = \prod_{i=1}^{n} q_{ic_i}, \sum_k q_{ik} = 1, 1 \le i \le n\}$. The subset Q contains all the distributions under which \mathbf{c} is jointly independent.

• Denote the maximizer in (7) by $(\check{B},\check{\beta})$ and \check{c} .

► SBM: (Daudin et al., 2008; Celisse et al., 2012; Bickel et al., 2013).

$$\max_{\boldsymbol{\theta}} \log P(\boldsymbol{A}, \boldsymbol{X}; \boldsymbol{\theta}) = \max_{\boldsymbol{\theta}, \boldsymbol{Q}(\cdot)} \mathbb{E}_{Q}[\log P(\boldsymbol{A}, \boldsymbol{X}, \mathbf{c}; \boldsymbol{\theta}) - \log Q(\mathbf{c})], \quad (6)$$

where $Q(\cdot)$ denotes any joint distribution on **c**.

► Here, we consider the mean-field variational approach (Jordan et al., 1999),

$$\max_{\boldsymbol{\theta}, \boldsymbol{Q} \in \boldsymbol{\mathcal{Q}}} \mathbb{E}_{Q}[\log P(A, X, \mathbf{c}; \boldsymbol{\theta}) - \log Q(\mathbf{c})],$$
(7)

where $Q = \{Q : Q(\mathbf{c}) = \prod_{i=1}^{n} q_{ic_i}, \sum_k q_{ik} = 1, 1 \le i \le n\}$. The subset Q contains all the distributions under which \mathbf{c} is jointly independent.

• Denote the maximizer in (7) by $(\check{B},\check{\beta})$ and \check{c} .

Theorem 2

Suppose the conditions of Theorem 1 hold. Then as $n o \infty$

$$P(\check{\boldsymbol{c}} = \boldsymbol{c}) \rightarrow 1, \quad \sqrt{n}(\check{\boldsymbol{\beta}} - \boldsymbol{\beta}) \stackrel{d}{\rightarrow} N(\boldsymbol{0}, I^{-1}(\boldsymbol{\beta})),$$

where $I(\beta)$ is the Fisher information for the multi-logistic regression problem of regressing c on X.

Remark: the same asymptotics as the MLE when $\frac{n\rho_n}{\log n} \to \infty$ as long as the network has sufficient edge information for doing approximate inference.

Idea: treat the community assignments \mathbf{c} as an unknown parameter to estimate.

$$(\tilde{\beta}, \tilde{B}, \tilde{\mathbf{c}}) = \underset{\substack{\beta_{K} = \mathbf{0}, \ \beta \in \mathbb{R}^{Kp} \\ \mathbf{c} \in \{1, \dots, K\}^{n} \\ B \in [0,1]^{K \times K}, B^{T} = B}}{\arg \max} \prod_{i < j} B_{c_{i}c_{j}}^{A_{ij}} (1 - B_{c_{i}c_{j}})^{1 - A_{ij}} \cdot \prod_{i} \frac{e^{\beta_{c_{i}}^{A} \mathbf{x}_{i}}}{\sum_{k=1}^{K} e^{\beta_{k}^{T} \mathbf{x}_{i}}}.$$
 (8)

Reference for SBM and DCBM: Bickel and Chen (2009); Zhao et al. (2012). For theoretical convenience, we consider a slightly different formulation:

$$(\tilde{\boldsymbol{\beta}}, \tilde{\boldsymbol{E}}, \tilde{\mathbf{c}}) = \underset{\substack{\boldsymbol{\beta_{K}=0, \ \boldsymbol{\beta}\in\mathbb{R}^{K_{p}}\\ \mathbf{c}\in\{1,\dots,K\}^{n}\\ \boldsymbol{\beta\in\mathbb{R}^{K\times K}, \boldsymbol{\beta}^{T}=\boldsymbol{B}}}}{\arg\max} \prod_{i< j} e^{-B_{c_{i}c_{j}}} B_{c_{i}c_{j}}^{A_{ij}}} \cdot \prod_{i} \frac{e^{\boldsymbol{\beta}_{c_{i}}^{T} \boldsymbol{x}_{i}}}{\sum_{k=1}^{K} e^{\boldsymbol{\beta}_{k}^{T} \boldsymbol{x}_{i}}}.$$
(9)

- T

$$\beta^{t+1} = \underset{\beta \in \mathbb{R}^{\rho K}, \beta_{K}=0}{\arg \max} \sum_{i} \left[\beta^{T}_{c_{i}^{t}} \mathbf{x}_{i} - \log\left(\sum_{k} e^{\beta^{L}_{k} \mathbf{x}_{i}}\right) \right],$$
(10)
$$B^{t+1}_{ab} = \underset{B_{ab}}{\arg \max} \log B_{ab} \cdot \sum_{i < j} A_{ij} \mathbf{1} (c_{i}^{t} = a, c_{j}^{t} = b) - B_{ab} \cdot \sum_{i < j} \mathbf{1} (c_{i}^{t} = a, c_{j}^{t} = b),$$
(11)
$$\mathbf{c}^{t+1} = \underset{\mathbf{c} \in \{1, \dots, K\}^{n}}{\arg \max} \sum_{ab} \left[\log B^{t+1}_{ab} \cdot \sum_{i < j} A_{ij} \mathbf{1} (c_{i} = a, c_{j} = b) - B_{ab} \cdot \sum_{i < j} \mathbf{1} (c_{i}^{t} = a, c_{j}^{t} = b),$$
(12)
$$- B^{t+1}_{ab} \cdot \sum_{i < j} \mathbf{1} (c_{i} = a, c_{j} = b) \right] + \sum_{i} (\beta^{t+1}_{c_{i}})^{T} \mathbf{x}_{i}.$$

Pairwise covariates-adjusted block model for community detection, IMS workshop at NUS $\triangleleft \square \flat \triangleleft \square \flat \triangleleft \blacksquare \flat \triangleleft \blacksquare \flat \triangleleft \blacksquare \flat \square \blacksquare \square \square \square \square \square$

Theorem 3

Assume the data (A, X) follows NSBM and Conditions 1 and 2 hold.

(i) If $n\rho_n \to \infty$ and $\mathbb{E} \| \mathbf{x} \|_2^{\alpha} < \infty$ ($\alpha > 1$), then there exists a constant $\gamma > 0$ such that, as $n \to \infty$

$$P\Big(\frac{1}{n}\sum_{i=1}^n \mathbb{1}(\tilde{c}_i \neq c_i) \leq \gamma(n\rho_n)^{-1/2}\Big) \to 1, \quad \|\tilde{\beta} - \beta\|_2 = O_p((n\rho_n)^{\frac{1-\alpha}{2\alpha}}).$$

(ii) Assume Condition 3 is satisfied. If $\frac{n\rho_n}{\log n} \to \infty$, then as $n \to \infty$

$$P(\tilde{\boldsymbol{c}} = \boldsymbol{c})
ightarrow 1, \quad \sqrt{n}(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}) \stackrel{d}{
ightarrow} N(\boldsymbol{0}, I^{-1}(\boldsymbol{\beta})),$$

where $I(\beta)$ is the Fisher information for the multi-logistic regression problem of regressing c on X.

 Non-convex. Multiple local optima may exist and the global solution is often impossible to accurately allocate.

1

- A good convex initialization is desirable.
- A popular class of convex initialization: semidefinite programming (SDP). Chen et al. (2012); Amini and Levina (2014); Cai and Li (2015); Montanari and Sen (2015); Guedon and Vershynin (2015).

$$\hat{Z} = \underset{Z}{\operatorname{arg\,max}} \quad \langle A + \gamma_n X X^T, Z \rangle \tag{13}$$
subject to
$$Z \succeq 0, Z \in \mathbb{R}^{n \times n}$$

$$0 \leq Z_{ij} \leq 1, 1 \leq i, j \leq n$$

$$\sum_{ij} Z_{ij} = \lambda_n,$$

where $\gamma_n, \lambda_n > 0$ are two tuning parameters; $\langle \cdot, \cdot \rangle$ denotes the inner product of two matrices. Similar spirit as Binkiewicz et al. (2014) for spectral clustering. We then obtain the communities by running K-means on \hat{Z} (treating each row of \hat{Z} as a data point in \mathbb{R}^n).

Solving (13) via alternating direction method of multipliers (ADMM) (Boyd et al., 2011).

Theory for SDP

Theorem 4

Assume condition (b) in NSBM holds and $(c_1, \mathbf{x}_1), \ldots, (c_n, \mathbf{x}_n) \stackrel{iid}{\sim} (c, \mathbf{x})$. Let $\{\bar{c}_i\}_{i=1}^n$ be the community estimates returned by running K-means on \hat{Z} . If $\min_a \bar{B}_{aa} > \max_{a \neq b} \bar{B}_{ab}, n\rho_n \to \infty$ and $\|\mathbf{x}\|_2$ is sub-Gaussian, then by choosing $\gamma_n = O(n^{-1})$ and $\lambda_n = \sum_{k=1}^K (\sum_{i=1}^n 1(c_i = k))^2$, we have

$$\frac{1}{n}\sum_{i=1}^n \mathbb{1}(\bar{c}_i\neq c_i)\stackrel{P}{\to} 0.$$

Remark:

- ► The crucial assumption min_a B_{aa} > max_{a≠b} B_{ab} requires denser edge connections within communities than between them.
- ▶ The tuning γ_n trades off the information from two different sources : network edge and nodal covariates. The specification $\gamma_n = O(n^{-1})$ in Theorem 4 implies that a larger weight should be given to network edge information.

Solving (13) via alternating direction method of multipliers (ADMM) (Boyd et al., 2011).

Input: initialize Z⁰ = A + γ_nXX^T, W⁰ = Y⁰ = U⁰ = V⁰ = 0, number of iterations *T*, step size ξ.

► For
$$t = 0, ..., T - 1$$

(a) $Y^{t+1} = \min\{\max\{0, \frac{1}{2}(W^t + Z^t - U^t - V^t)\}, 1\}$
(b) $W^{t+1} = Y^{t+1} + U^t + \frac{\lambda_n - \sum_{ij}(Y_{ij}^{t+1} + U_{ij}^t)}{n^2} \mathbf{11}^T$
(c) $Z^{t+1} = P\Lambda_+P^T$, where $Y^{t+1} + V^t + \xi^{-1} \cdot (A + \gamma_n XX^T) = P\Lambda P^T$
(d) $U^{t+1} = U^t + Y^{t+1} - W^{t+1}, V^{t+1} = V^t + Y^{t+1} - Z^{t+1}$
► Output Z^{T+1}