

Achieving Optimal Misclassification Proportion in Stochastic Block Model - Review

Author: Chao Gao, Zongming Ma, Anderson Y.Zhang, Harrison H.Zhou

Jiixin Hu

First edition: 04/19/2020; Final edition: 04/22/2020

ABSTRACT

This paper proposes a polynomial time two-stage method to detect the community in network matrix data with stochastic block model(SBM). One stage is to initialize the partition with the provided new greedy clustering method or any other methods with weak convergence condition. The other stage is to refining the node-wise partition via penalized local maximum likelihood estimation. Statistical guarantees for refinement scheme and initialization are given. The paper reviews related literature closely, which provides a good view of network community detection.

1 PROBLEM FORMULATION & MODEL

1.1 Model

Let $A \in \{0,1\}^{n \times n}$ be the symmetric adjacency matrix where $A_{ii} = 0, i \in [n]$ and $A_{uv} = A_{vu}, \forall u > v \in [n]$. Assume there are k communities among n nodes. The stochastic block model can be written as below:

$$A_{uv} \sim_{i.i.d} \text{Ber}(P_{uv}), \forall u \neq v; \quad \mathbb{E}(A_{uv}) = P_{uv} = B_{\sigma(u)\sigma(v)},$$

where $B \in [0,1]^{k \times k}$ is the connectivity matrix and $\sigma : [n] \rightarrow [k]$ is the label function. In matrix form, the model can also be presented as below:

$$\mathbb{E}(A) = P = H^T B H; \quad A = H^T B H + \mathcal{E},$$

where $H \in \{0,1\}^{k \times n}$ is the membership matrix corresponding to the partition σ and \mathcal{E} is a sub-Gaussian noise matrix. This model can be considered as the order-2 case of TBM(M.Wang 2019) and the single layer case of TSBM(J.Lei 2020).

Two parameter spaces are discussed in the paper:

$$\Theta_0(n, k, a, b, \beta) = \left\{ (B, \sigma) \mid \sigma : [n] \rightarrow [k], n_i = \sum_{u \in [n]} \mathbb{I}\{\sigma(u) = i\} \in \left[\frac{n}{\beta k} - 1, \frac{\beta n}{k} + 1 \right], \forall i \in [k], \right. \\ \left. B = [B_{ij}] \in [0,1]^{k \times k}, B_{ii} = \frac{a}{n} \text{ for } \forall i, B_{ij} = \frac{b}{n} \text{ for } \forall i \neq j \right\},$$

where $\beta \geq 1$ is an absolute constant which controls the range of community size. Relaxing the equal

within and equal between connection, we have a larger parameter space:

$$\begin{aligned}\Theta(n, k, a, b, \lambda, \beta, \alpha) = & \left\{ (B, \sigma) \mid \sigma : [n] \rightarrow [k], n_i = \sum_{u \in [n]} \mathbb{I}\{\sigma(u) = i\} \in \left[\frac{n}{\beta k} - 1, \frac{\beta n}{k} + 1 \right], \forall i \in [k], \right. \\ & B = B^T = \llbracket B_{ij} \rrbracket \in [0, 1]^{k \times k}, \frac{b}{\alpha n} \leq \frac{1}{k(k-1)} \sum_{i \neq j} B_{ij} \leq \max_{i \neq j} B_{ij} = \frac{b}{n}, \\ & \left. \frac{a}{n} = \min_i B_{ii} \leq \max_i B_{ii} \leq \frac{\alpha a}{n}, \lambda_k(P) \leq \lambda \text{ with } P = \llbracket P_{uv} \rrbracket = \llbracket B_{\sigma(u)\sigma(v)} \rrbracket \right\},\end{aligned}$$

where α is absolute constant which controls the range of signal.

Additionally, authors assume $0 < \frac{b}{n} < \frac{a}{n} \leq 1 - \epsilon$ for some constant $\epsilon \in (0, 1)$ throughout the paper. That means $0 < \max_{i \neq j} B_{ij} < \min_i B_{ii} < 1 - \epsilon$. To ensure $\Theta_0(n, k, a, b, \beta) \subset \Theta(n, k, a, b, \lambda, \beta, \alpha)$, the paper also requires $\lambda \leq \frac{a-b}{2\beta k}$ throughout the paper.

This main goal for this paper is to find an estimate of partition $\hat{\sigma}$ to achieve the optimal minimax misclassification proportion established in Zhang and Zhou(2015).

1.2 Discussion of Loss measure

First, the loss measure in Gao's and Zhang's paper is stated below:

$$l(\hat{\sigma}, \sigma) = \min_{\pi \in S_k} \frac{1}{n} \sum_{u \in [n]} \mathbb{I}\{\hat{\sigma}(u) \neq \pi(\sigma(u))\},$$

where S_k stands for the symmetric group on $[k]$ consisting of all permutations of $[k]$ and $\pi(\sigma(\cdot))$ refers a permuted σ . Recall the misclassification proportion in TBM and TBSM:

$$MCR(\hat{M}_k, M_{k, \text{true}}) = \max_{r \in [R_k], a \neq a' \in [R_k]} \min\{D_{a,r}^{(k)}, D_{a',r}^{(k)}\}, \text{ where } D_{r,r'}^{(k)} = \frac{1}{d_k} \sum_{i=1}^{d_k} \mathbb{I}[m_{i,r}^{(k)} = \hat{m}_{i,r'}^{(k)} = 1], r, r' \in [R_k].$$

To show the relationship between $l(\hat{\sigma}, \sigma)$ and $MCR(\hat{M}_k, M_{k, \text{true}})$, define:

$$D_{r,r'} = \sum_{u \in [n]} \mathbb{I}\{\hat{H}_{ur} = H_{ur'} = 1\}; \quad D_{r,r'}^\pi = \sum_{u \in [n]} \mathbb{I}\{\hat{H}_{ur} = \pi(H_{ur'}) = 1\}$$

For here is an order-2 symmetric case, there is only one confusion matrix D for TBM. Therefore, we can show

$$MCR(\hat{H}, H) = \frac{1}{n} \max_{a \neq a' \in [k], r \in [k]} \min\{D_{ar}, D_{a'r}\}.$$

Let $n_r = \sum_{u \in [n]} \mathbb{I}\{\pi(\sigma(u)) = r\}$, $r \in [k]$ and $\pi(H)$ refers to the membership matrix corresponds to $\pi(\sigma)$. Then we have,

$$\begin{aligned}
\sum_{u \in [n]} \mathbb{I}\{\hat{\sigma}(u) \neq \pi(\sigma(u))\} &= \sum_{n \in [n]} \sum_{r=1}^k \mathbb{I}\{\hat{H}_{ur} = 0, \pi(H_{ur}) = 1\} \\
&= \sum_{r=1}^k \sum_{n \in [n]} \mathbb{I}\{\hat{H}_{ur} = 0, \pi(H_{ur}) = 1\} \\
\text{Because } \sum_{u \in [n]} \mathbb{I}\{\pi(H_{ur}) = 1\} &= n_r, = \sum_{r=1}^k \left(n_r - \sum_{u \in [n]} \mathbb{I}\{\hat{H}_{ur} = 1 = \pi(H_{ur})\} \right) \\
&= \sum_{r=1}^k (n_r - D_{r,r}^\pi) \\
&\geq^* \sum_{r=1}^k (n_r - \max_{r' \in [k]} D_{r',r}) \\
&\geq \sum_{r=1}^k (\max_{a, a' \in [k]} \min\{D_{ar}, D_{a'r}\}) \\
&\geq \max_{a, a' \in [k], r \in [k]} \min\{D_{ar}, D_{a'r}\},
\end{aligned}$$

where the inequality \geq^* is not necessarily become equality even though $\pi = \pi^0 = \arg \min_{\pi \in S_k} \sum_{u \in [n]} \mathbb{I}\{\hat{\sigma}(u) \neq \pi(\sigma(u))\}$. Take a toy example. Let $k = 3$ and suppose below confusion matrix:

$$D = \begin{bmatrix} 10 & 9 & 8 \\ 5 & 6 & 6 \\ 5 & 5 & 6 \end{bmatrix}.$$

Here $D^{\pi^0} = D$ while $D_{r,r}^{\pi^0} < \max_{r' \in [3]} D_{r',r}$ for $r = 2, 3$. We can also calculate $MCR(\hat{H}, H) = \frac{6}{60}$ while $l(\hat{\sigma}, \sigma) = \frac{5+5+9+5+8+6}{60}$. We can conclude that $l(\hat{\sigma}, \sigma) \geq MCR(\hat{H}, H)$. And when the clusters k is fixed, there are $l(\hat{\sigma}, \sigma) \asymp MCR(\hat{H}, H)$. When $k \rightarrow \infty$, for $l(\hat{\sigma}, \sigma)$:

$$l(\hat{\sigma}, \sigma) = \frac{\sum_{r \neq r', r, r' \in [k]} D_{r,r'}^\pi}{n} = \frac{k(k-1)\bar{D}_{r,r'}^\pi}{n},$$

where $\bar{D}_{r,r'}^\pi$ is the average of $D_{r,r'}^\pi$, $\forall r \neq r' \in [k]$, which has value at the same level of the entry in D . Here I am wonder how does the entry of D change along with k, n . Consider the worst case of random guess when the community sizes are nearly equal, $D_{r,r'} \asymp \frac{n}{k^2}$. From this perspective, $MCR \asymp \frac{D_{r,r'}}{n} \asymp \frac{1}{k^2}$ will goes to infinity as long as $k \rightarrow \infty$ while $l(\hat{\sigma}, \sigma) = O(1)$ when $k \rightarrow \infty, n \rightarrow \infty$ for the worst case.

Similarly, for the $\eta = MCR = \frac{1}{n_{min}} \max_{a \neq a' \in [k], r \in [k]} \min\{D_{ar}, D_{a'r}\}$ in Lei's paper, consider the nearly equal community sizes case, $n_{min} \asymp \frac{n}{k}$. We still have $\eta \asymp \frac{D_{r,r'}^k}{n} \asymp \frac{1}{k}$. That also implies η in Lei's paper will vanish when $k \rightarrow \infty$ even though the worst case.

1.3 Discussion of Optimal Misclassification proportion

The optimal minimax misclassification proportion is established by Zhang and Zhou(2015).

To show the optimal rate, define:

$$I^* = -2\log\left(\sqrt{\frac{a}{n}}\sqrt{\frac{b}{n}} + \sqrt{1-\frac{a}{n}}\sqrt{1-\frac{b}{n}}\right),$$

which is exactly the Renyi Divergence $D_{1/2}(Ber(a/n)\|Ber(b/n))$. The definition of order 1/2 Renyi Divergence is $D_{1/2}(P\|Q) = -2\log(\sum_i^n \sqrt{p_i q_i})$ where $i \in \mathbb{Z}$ is the outcome of discrete random variable X and $P(X = i) = p_i$. According to the Lemma B.1 in the supplement of Zhang and Zhou(2015), $I^* \asymp \frac{(a-b)^2}{na}$. That makes sense for the distance between the two Bernoulli distributions decreases when $n \rightarrow \infty$ and $a = O(1)$.

Therefore the optimal minimax misclassification rate is showed in below theorem.

Theorem 1 (Optimal rate for $\hat{\sigma}$, Zhang and Zhou(2015)). *Assume $\frac{nI^*}{k \log k} = \frac{(a-b)^2}{ak \log k} \rightarrow \infty$, then*

$$\inf_{\hat{\sigma}} \sup_{\Theta} \mathbb{E}l(\sigma, \hat{\sigma}) = \begin{cases} \exp\left(-(1+o(1))\frac{nI^*}{2}\right), & k = 2 \\ \exp\left(-(1+o(1))\frac{nI^*}{\beta k}\right), & k \geq 3 \end{cases},$$

for both $\Theta = \Theta_0(n, k, a, b, \beta)$ and $\Theta = \Theta(n, k, a, b, \lambda, \beta, \alpha)$ with $\lambda \leq \frac{a-b}{2\beta k}$, where $\beta \in [1, \sqrt{5/3}]$.

Note that **Thm 1** requires $\frac{(a-b)^2}{ak \log k} \rightarrow \infty$. That means if a is a constant, the connectivity signal is $O(n^{-1})$ and $nI^* = O(1)$, where no consistent estimate $\hat{\sigma}$ exists. Therefore, to have a consistent estimate, at least $a > O(1)$.

When k is fixed, the $l(\hat{\sigma}, \sigma)$ (called MCR_{opt} below) converges to 0 as $\exp(-a) \rightarrow 0$, which is equal to $\exp(a) \rightarrow \infty$. Recall the asymptotic result in TBM and TBSM, the MCR_{TBM} and MCR_{TBSM} goes to 0 as $n \rightarrow \infty$. If $O(1) < a < O(\log n)$, $\exp(a)$ is slower than n to go to the infinity. If $a = O(\log n)$, MCR_{opt} has the same rate as $MCR_{TBM, TBSM}$. If $a > O(\log n)$, then MCR_{opt} is faster than the other two. Notice that in TBM and TBSM, we require the irreducibility of the connectivity, however, in Θ_0 and Θ , if $a < O(n)$, then the signal a/n and b/n will vanish. To also satisfies the irreducibility, we should require $a \geq O(n)$. Therefore, when 3 models all have irreducibility, MCR_{opt} is much faster than $MCR_{TBM, TBSM}$.

When k is divergent, MCR_{TBM} and MCR_{TBSM} are unstable and degenerate, we can not tell any conclusion about them this time. For MCR_{opt} , the model allows $k \rightarrow \infty$ if satisfies $\frac{(a-b)^2}{ak \log k} \rightarrow \infty$. Therefore, if $k \log k < O(a)$, the misclassification proportion will converge to 0 as $\exp(a/k) \rightarrow \infty$.

According to Zhang and Zhou(2015), the optimal lower bound is achieved by a novel reduction of the *global* minimax rate in to a *local* testing problem. Different with the regular meaning of *global* in *global optimum*, here *global* refers to that the optimal rate is a *global* property for the **whole** network. In contrast, *local* loss will only focuses on one node, i.e. node-wise loss. Lemma 2.1 in Zhang and Zhou(2015) makes it possible to study the global optimum via local loss when the model is homogeneous ($B_{ii} = a/n, B_{i,j} = b/n, \forall i \neq j \in [k]$) and close under permutation.

Lemma 2 (Global to Local). *Let Γ be any parameter space that SBM is homogeneous and close under permutation. Let $S_{\sigma}(\hat{\sigma}) = \{\sigma' : \sigma' = \pi \circ \hat{\sigma}, \pi = \arg \min_{\pi \in S_k} \sum_{u \in [n]} \mathbb{I}\{\pi(\hat{\sigma}(u)) \neq \sigma(u)\}\}$ and $l(\hat{\sigma}(i), \sigma(i)) =$*

$\sum_{\sigma' \in S_\sigma(\hat{\sigma})} \frac{\mathbb{I}\{\sigma'(i) \neq \sigma(i)\}}{|S_\sigma(\hat{\sigma})|}$. Assume $L(\hat{\sigma}) = \frac{1}{|\Gamma|} \sum_{\sigma \in \Gamma} \mathbb{E}l(\sigma, \hat{\sigma})$ and $L(\hat{\sigma}(1)) = \frac{1}{|\Gamma|} \sum_{\sigma \in \Gamma} \mathbb{E}l(\sigma(1), \hat{\sigma}(1))$. Then we have:

$$\inf_{\hat{\sigma}} L(\hat{\sigma}) = \inf_{\hat{\sigma}} L(\hat{\sigma}(1))$$

Intuitively, for heterogeneous SBM, we can add penalty to find an optimal $\hat{\sigma}$ that also reaches the optimal rate. Indeed, according to Zhang and Zhou(2015), a range of penalized likelihood-type estimates can achieve the optimal rate. Inspired by the penalized likelihood method and the thoughts of global to local, this paper proposes a computational feasible algorithm whose estimates can be rate-optimal.

2 ALGORITHM

This paper proposes a two-stage algorithm which contains a initialization step and a refinement step. The corresponding asymptotic results are provided.

2.1 Procedures

2.1.1 Initialization

Though the theoretical results show that refinement stage only require the weak consistency condition of initialization method, authors here introduce a greedy spectral clustering algorithm.

First, consider the trimming unnormalized spectral clustering (USC) with adjacency matrix A . Let $d_u = \sum_{v \in [n]} A_{uv}$ be the degree of u -th node. Obtain the trimmed adjacency matrix $T_\tau(A)$ be replacing $A_{u\cdot}$ and $A_{\cdot u}$ to 0 if $d_u \geq \tau$.

Second, consider the trimming normalized spectral clustering (NSC) with graph Laplacian $L(A)$, where $\llbracket L(A)_{uv} \rrbracket = d_u^{-1/2} d_v^{-1/2} A_{uv}$. Obtain the trimmed Laplacian $L(A_\tau)$ by replacing A to $A_\tau = A + \frac{\tau}{n} \mathbf{1}\mathbf{1}^T$.

The sketch of the greedy clustering method (Algorithm 1) is below:

Algorithm 1 Initialization

Input: Data matrix $\hat{U} = T_\tau(A)$ or $L(A_\tau)$, # of communities k , critical radius $r = \mu \sqrt{\frac{k}{n}}$ with some constant μ .

Output: Community assignment $\hat{\sigma}$.

- 1: $S = [n]$;
 - 2: **for** $i = 1$ to k **do**
 - 3: $t_i = \arg \max_{u \in S} \sum_{v \in S} \mathbb{I}\{\|\hat{U}_v - \hat{U}_u\| \leq r\}$;
 - 4: $\hat{\mathcal{C}}_i = \{v \in S : \|\hat{U}_v - \hat{U}_{t_i}\| \leq r\}$;
 - 5: Label $\hat{\sigma}(u) = i$ for all $u \in \hat{\mathcal{C}}_i$;
 - 6: $S \leftarrow S \setminus \hat{\mathcal{C}}_i$.
 - 7: **end for**
 - 8: If $S \neq \emptyset$, then for any $u \in S$, $\hat{\sigma}(u) = \arg \min_{i \in [k]} \frac{1}{|\hat{\mathcal{C}}_i|} \sum_{v \in \hat{\mathcal{C}}_i} \|\hat{U}_u - \hat{U}_v\|$.
-

2.1.2 Refinement

Under the parameter space $\Theta_0(n, k, a, b, 1)$ with equal community size, the MLE for σ is:

$$\hat{\sigma} = \arg \max_{\sigma} \sum_{u < v} A_{uv} \mathbb{I}\{\sigma(u) = \sigma(v)\},$$

which is a combinatorial optimization problem and is computationally intractable.

However, we can easily write the close form solution for the node-wise optimization. If we know the values of $\{\sigma(u)\}_{u=2}^n$, then for $\sigma(1)$:

$$\hat{\sigma}(1) = \arg \max_{i \in [k]} \sum_{v \neq 1: \sigma(v)=i} A_{1v}.$$

The estimate can be interpreted as the first node should belong to the community where the first node has the most neighbours. In practice, we do know the labels in advance. However, we can estimate the labels for $(n-1)$ nodes with initialization algorithm. We note σ^0 as initialization algorithm for the $A_{-u} \in \{0, 1\}^{(n-1) \times (n-1)}$, where A_{-u} is a submatrix of A with its u -th row and columns are removed. Repeat this node-wise optimization among every node, we get the refinement scheme for community detection (Algorithm 2).

Algorithm 2 Refinement

Input: Adjacency matrix A , # of communities k , initialization algorithm σ^0 .

Output: Community assignment $\hat{\sigma}$.

1: **for** $u = 1$ to n **do**

2: Apply σ^0 to A_{-u} , obtain $\sigma_u^0(v)$, $\forall v \neq u$ and let $\sigma_u^0(u) = 0$

3: Define $\tilde{\mathcal{C}}_i^u = \{v : \sigma_u^0(v) = i\}$, let $\tilde{\mathcal{E}}_i^u = \sum_{x < y, \sigma_u^0(x), \sigma_u^0(y) \in \tilde{\mathcal{C}}_i^u} A_{xy}$ and $\tilde{\mathcal{E}}_{ij}^u = \sum_{\sigma_u^0(x) \in \tilde{\mathcal{C}}_i^u, \sigma_u^0(y) \in \tilde{\mathcal{C}}_j^u} A_{xy}$.

4: Define

$$\hat{B}_{ii}^u = \frac{|\tilde{\mathcal{E}}_i^u|}{\frac{1}{2} |\tilde{\mathcal{C}}_i^u| (|\tilde{\mathcal{C}}_i^u| - 1)}, \quad \hat{B}_{ij}^u = \frac{|\tilde{\mathcal{E}}_{ij}^u|}{|\tilde{\mathcal{C}}_i^u| |\tilde{\mathcal{C}}_j^u|}, \quad \forall i \neq j \in [k]$$

5: Define $\hat{\sigma}_u(v) = \sigma_u^0(v)$ for all $i \neq j$, and

$$\hat{\sigma}_u(u) = \arg \max_{l \in [k]} \sum_{\sigma_u^0(v)=l} A_{uv} - \rho_u \sum_{v \in [n]} \mathbf{1}_{\{\sigma_u^0(v)=l\}},$$

where

$$t_u = \frac{1}{2} \log \frac{\hat{a}_u (1 - \hat{b}_u/n)}{\hat{b}_u (1 - \hat{a}_u/n)}, \quad \rho_u = -\frac{1}{2t_u} \log \left(\frac{\frac{\hat{a}_u}{n} e^{-t_u} + 1 - \frac{\hat{a}_u}{n}}{\frac{\hat{b}_u}{n} e^{t_u} + 1 - \frac{\hat{b}_u}{n}} \right).$$

6: **end for**

7: **Consensus:** Define $\hat{\sigma}(1) = \hat{\sigma}_1(1)$, for $u = 2, \dots, n$, define

$$\hat{\sigma}(u) = \arg \max_{l \in [k]} |\{v : \hat{\sigma}_1(v) = l\} \cap \{v : \hat{\sigma}_u(v) = \hat{\sigma}_u(u)\}|.$$

2.2 Asymptotic Result for Algorithm

2.2.1 Theorem for Initialization

Theorem 3 (Performance for Unnormalized Spectral Clustering). Assume $e \leq a \leq C_1 b$ for some constant $C_1 > 0$ and $\frac{ka}{\lambda_k^2} \leq c$ for sufficient small $c \in (0, 1)$. Consider $USC(\tau)$ with sufficient small μ and $\tau = C_2 \bar{d}$, where $\bar{d} = 1/n \sum_{u \in [n]} d_u$ and sufficient large constant $C_2 > 0$. For any C' , there exists some $C > 0$ depends on C_1, C_2, C', μ s.t.

$$\ell(\widehat{\sigma}, \sigma) \leq C \frac{a}{\lambda_k^2}$$

with probability at least $1 - n^{-C'}$. If k is fixed, the same conclusion holds without $a \leq C_1 b$.

Theorem 4 (Performance for Normalized Spectral Clustering). Assume $e \leq a \leq C_1 b$ for some constant $C_1 > 0$ and $\frac{ka \log(a)}{\lambda_k^2} \leq c$ for sufficient small $c \in (0, 1)$. Consider $USC(\tau)$ with sufficient small μ and $\tau = C_2 \bar{d}$, where $\bar{d} = 1/n \sum_{u \in [n]} d_u$ and sufficient large constant $C_2 > 0$. For any C' , there exists some $C > 0$ depends on C_1, C_2, C', μ s.t.

$$\ell(\widehat{\sigma}, \sigma) \leq C \frac{a \log(a)}{\lambda_k^2}$$

with probability at least $1 - n^{-C'}$. If k is fixed, the same conclusion holds without $a \leq C_1 b$.

These two theorem are suitable for both Θ and Θ_0 .

2.2.2 Theorem for Refinement

Condition 2.1. There exists constants $C_0, \delta > 0$ and positive sequence $\gamma = \gamma_n$, s.t.

$$\inf_{(B, \sigma) \in \Theta} \min_{u \in [n]} \mathbb{P}_{B, \sigma} \left\{ \ell(\sigma, \sigma_u^0) \leq \gamma \right\} \geq 1 - C_0 n^{-(1+\delta)}$$

for some Θ .

Theorem 5 (Performance for refinement). Suppose as $n \rightarrow \infty$, $\frac{(a-b)^2}{ak \log k} \rightarrow \infty$, $a \asymp b$ and Condition 1 satisfies for **1**: $\gamma = o\left(\frac{1}{k \log k}\right)$ and $\Theta = \Theta_0(n, k, a, b, \beta)$ **2**: $\gamma = o\left(\frac{1}{k \log k}\right)$ and $o\left(\frac{a-b}{ak}\right)$ and $\Theta = \Theta(n, k, a, b, \lambda, \beta, \alpha)$. Then there is a sequence $\eta \rightarrow 0$ s.t.

$$\sup_{(B, \sigma) \in \Theta} \mathbb{P}_{B, \sigma} \left\{ \ell(\sigma, \widehat{\sigma}) \geq \exp\left(- (1 - \eta) \frac{nI^*}{2}\right) \right\} \rightarrow 0, \quad \text{if } k = 2$$

$$\sup_{(B, \sigma) \in \Theta} \mathbb{P}_{B, \sigma} \left\{ \ell(\sigma, \widehat{\sigma}) \geq \exp\left(- (1 - \eta) \frac{nI^*}{\beta k}\right) \right\} \rightarrow 0, \quad \text{if } k \geq 3,$$

which can be rewritten as :

$$\ell(\sigma, \widehat{\sigma}) \leq \exp\left(- (1 - \eta) \frac{nI^*}{2}\right), \quad \text{if } k = 2$$

$$\ell(\sigma, \widehat{\sigma}) \leq \exp\left(- (1 - \eta) \frac{nI^*}{\beta k}\right), \quad \text{if } k \geq 3$$

with probability goes to 1 as $n \rightarrow \infty$, $\frac{(a-b)^2}{ak \log k} \rightarrow \infty$.

2.3 Discussion of the Theorems

When $k = O(1)$, $\gamma = o(\frac{1}{k \log k})$ reduces to $\gamma = o(1)$ and $\gamma = \frac{a-b}{ak}$ reduces to $\gamma = o(\frac{a-b}{a}) = o(1)$. For Θ_0 , we have $\lambda_k \geq \frac{a-b}{\beta k}$. To fulfill **Thm 3**'s condition $\frac{ka}{\lambda^2} \leq c$, we only requires $\frac{(a-b)^2}{a} \rightarrow \infty$ and then satisfies **1** for **Thm 5**. That implies when $k = O(1)$, $\frac{(a-b)^2}{a} \rightarrow \infty$, the algorithm initializes with USC and refines as Algorithm 2 will reach the optimal minimax misclassification rate. Similarly, if $\frac{(a-b)^2}{a \log a} \rightarrow \infty$, **Thm 4** can be applied and NSC can also be a good initialization for the rate-optimal algorithm. And for Θ , $\lambda_k \leq \frac{a-b}{2\beta k}$. As above discussion, **2** in **Thm 5** can be satisfied if $\frac{(a-b)^2}{a} \rightarrow \infty$ for USC and $\frac{(a-b)^2}{a \log a} \rightarrow \infty$ for NSC.

3 TO DO LIST & QUESTION

3.1 To do list

- Maybe can extend this work to higher-order case.
- Figure out why the theorem allows $k \rightarrow \infty$? Can we use similar idea in TBM and TBSM?

3.2 Questions

- As we need $a \rightarrow \infty$ to get a consistent estimate, can the $a > O(n)$? In that case, this model not only allows the number of clusters k goes to ∞ but also allows the largest entry of the connection goes to infinity.

Possible Answer: I think a can not goes faster than n , because any entry in B should be smaller than 1. Therefore, as TBM and TBSM, the signal level of B is upper bounded by 1.

- The meaning of ρ_u ?

Possible Answer: Maybe needs to read relative parts in Zhang and Zhou(2015).

- Where is α in the theorem?

Possible Answer: In Zhang and Zhou, the heterogeneous setting do not uses α . Need to investigate more where these is an α in this paper.