

Abstract

**Community Detection:
Fundamental Limits, Methodology, and Variational Inference**

Ye Zhang

2018

Network analysis has become one of the most active research areas over the past few years. A core problem in network analysis is community detection. In this thesis, we investigate it under Stochastic Block Model and Degree-corrected Block Model from three different perspectives: 1) the minimax rates of community detection problem, 2) rate-optimal and computationally feasible algorithms, and 3) computational and theoretical guarantees of variational inference for community detection.

**Community Detection:
Fundamental Limits, Methodology, and
Variational Inference**

A Dissertation
Presented to the Faculty of the Graduate School
of
Yale University
in Candidacy for the Degree of
Doctor of Philosophy

by
Ye Zhang

Dissertation Director: Harrison H. Zhou

May 2018

Copyright © 2018 by Ye Zhang
All rights reserved.

Contents

Acknowledgements	vii
1 Introduction	1
2 Stochastic Block Model	4
3 Minimax Rates	6
3.1 Hypothesis Testing for One Single Node	8
3.2 Minimax Lower Bound	10
3.3 Minimax Upper Bound	13
3.4 Extension	16
3.5 Proof of Propositions	17
4 Methodology	20
4.1 Hypothesis Testing for One Single Node: Revisit	20
4.2 Spectral Clustering	22
4.3 Rate-optimal and Computationally Feasible Algorithm	26
4.4 Extension	31
5 Variational Inference	33
5.1 Mean Field Variational Inference	34
5.2 A Bayesian Framework	35
5.3 Mean Field Approximation	36
5.3.1 Coordinate Ascent Variational Inference	38

5.3.2	Batch Coordinate Ascent Variational Inference	39
5.4	Theoretical Justifications	41
5.4.1	Loss Function	41
5.4.2	Ground Truth	42
5.4.3	Guarantees	42
5.5	Proofs	45
5.5.1	Proof of Theorem 5.1	45
5.5.2	Proof of Theorem 5.2	47
5.5.3	Proof of Theorem 5.3	49
5.5.4	Additional Lemmas and Propositions and Their Proofs	60
6	Generalization: Degree-corrected Block Model	68
6.1	Model	69
6.2	Minimax Risks	71
6.3	An Adaptive and Computationally Feasible Procedure	74
6.3.1	A Two-Stage Algorithm	74
6.3.2	Performance Guarantees	75

List of Figures

1.1	Visualizations of political blogosphere dataset [2]. <i>Left</i> : connections among blogs, which are colored blue (democratic) or red (republican) according to their political views. <i>Middle</i> : adjacency matrix where nodes are randomly ordered. <i>Right</i> : adjacency matrix where nodes are ordered by their political views. The edges within democratic blogosphere are colored blue while those within republican blogosphere are colored red.	2
5.1	Graphical model presentations of full Bayesian inference (<i>left</i> panel) and the mean field approximation (<i>right</i> panel) for community detection. The edges show the dependence among variables.	37

*Dedicated to my parents Chaohua and Guimei
for their endless support and love.*

Acknowledgements

Having Harrison Zhou as my adviser is one of the best things ever happened to me. During my Ph.D. years, I received tremendous care, inspiration, and support from him. I benefited immensely from his enthusiasm for research, his generosity of sharing ideas and thoughts, and his remarkable intuition in analyzing statistical problems. His advice is truly helpful to my graduate studies.

I cannot thank David Pollard enough. Presenting in his YPNG seminars and meeting him regularly to talk about my research were a great privilege to me. His guidance and feedback are extremely beneficial. He also demonstrated me how to become a great teacher.

I am so fortunate to have Zongming Ma as my collaborator. Working with him is an enormously rewarding experience. I benefited tremendously from his passion and dedication to research, and his statistical and algorithmic insights on projects.

I appreciate the suggestions from Daniel Spielman and Yihong Wu on my presentations. I thank Andrew Barron, Joseph Chang, Jay Emerson, Sekhar Tatikonda, and Sahand Negahban for providing fantastic courses. I am also grateful to Joann DeVecchio, JoAnn Falato, Elizavette Torres, and Karen Kavanaugh for their support.

I am so happy to have Chao Gao and Yu Lu around during my Ph.D. years. Working together and playing video games together on the third floor of 24 Hillhouse Ave are the happiest moments I had in the past years. Their philosophy about research and life keeps influencing me in a very positive way. I also thank Natalie Doss, Derek Feng, Tal Sarig, Mengjie Chen and Zhao Ren for helpful discussions on research.

Finally, I would like to thank my parents Chaohua and Guimei for their endless love and support.

Chapter 1

Introduction

Network science [15, 36, 40, 49] has become one of the most active research areas over the past few years. It has applications in many disciplines, for example, physics [41], sociology [51], biology [6], and Internet [4]. The observed networks can often be modeled as an instance of a random graph and the goal is to infer structures of the underlying generating process. A structure of particular interest is *community*: there is a partition of the graph nodes in some suitable sense so that each node belongs to a community. The so-called *community detection* is to recover such community structure from the observed networks.

For the purpose of illustration, we use a well-known example of political blogosphere dataset [2] collected before 2014 United State presidential election. Blogs are labeled democratic or republican according to their political views. Two blogs are connected if there exists a hyperlink between them. The right panel of Figure 1.1 shows that there exist far more connections within the democratic blogosphere and the republican blogosphere while the connections between these two are fewer. The aim of community detection is to infer the community structure from the adjacency matrix (e.g., the middle panel of Figure 1.1), as usually there is no additional covariate information and the nodes are arbitrarily ordered.

This thesis covers different aspects of the community detection problem, including *models*, *fundamental limits*, *methodology*, and *variational inference*, as follows.

- **Models.** One of the most important open problems in community detection is about modeling. Stochastic Block Model (SBM) [29] has been the most popular and most

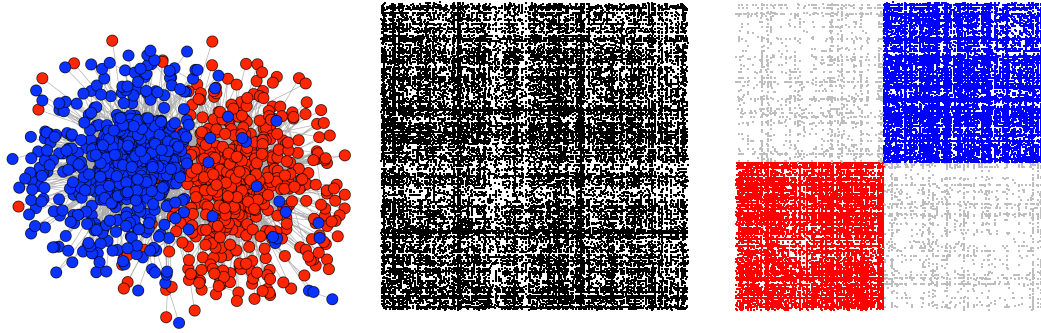


Figure 1.1: Visualizations of political blogosphere dataset [2]. *Left*: connections among blogs, which are colored blue (democratic) or red (republican) according to their political views. *Middle*: adjacency matrix where nodes are randomly ordered. *Right*: adjacency matrix where nodes are ordered by their political views. The edges within democratic blogosphere are colored blue while those within republican blogosphere are colored red.

well-investigated model in literature. It embeds communities on Erdős-Rényi graphs. It captures the community structure in a straightforward and intuitive way, and is simple enough to derive decent and beautiful theoretical results. This makes the SBM appealing to computer scientists, mathematicians, and statisticians. As a consequence, we mainly study the SBM in this thesis. We introduce the SBM in detail in Chapter 2, with further theoretical investigations presented in Chapters 3, 4, and 5. In addition, we extend all the results to a more general model, Degree-corrected Block Model (DCBM) [14, 33], in Chapter 6.

- **Fundamental Limits.** From a decision-theoretical point of view, once a model (SBM) is provided, the follow-up step is to derive minimax rates. In Chapter 3, we will show minimax rates of misclassification proportion takes an exponential form of signal-to-noise ratio, under minimum assumptions of parameters. On top of that, it immediately leads to various phase transitions and tight thresholds established in the literature [1, 10, 38].
- **Methodology.** Many methods have been implemented in practice for community detection problem, with spectral clustering, modularity, semi-definite programming being the most popular choices. In a general sense, there are two desired properties of an algorithm: 1) *computability*: whether it is a polynomial-time algorithm or

not; 2) *optimality*: whether it has provable theoretical guarantee that matches with minimax rate. Despite their popularity, none of the aforementioned methods satisfy both properties. In Chapter 4, we propose a novel two-stage algorithm that is both computationally feasible and rate-optimal. Our methods work for both sparse and dense networks.

- **Variational Inference.** Variational inference has been widely used to approximate posterior distributions. Despite popularity, it has very little theoretical justification established in the literature. In Chapter 5, we provide statistical and computational guarantees of variational inference for the SBM.

This thesis incorporates main results from our papers [19, 20, 52, 53]. Roughly speaking, there is a one-to-one correspondence between the chapters in this thesis and our papers: Chapter 3 is for [52], Chapter 4 is for [19], Chapter 5 is for [20], and Chapter 6 is for [53]. However, in order to capture the main ideas and contributions of our work, we will make some simplification to the the results in the aforementioned paper. In this way, we are able to better present our result without being overwhelmed by technical details.

Chapter 2

Stochastic Block Model

The SBM, proposed by [29], is the most studied model in community detection. Consider an n -node network with its adjacency matrix denoted by A . It is an unweighted and undirected network without self-loops, with $A \in \{0, 1\}^{n \times n}$, $A = A^T$ and $A_{i,i} = 0, \forall i \in [n]$. Each edge is an independent Bernoulli random variable with $\mathbb{E}A_{i,j} = P_{i,j}, \forall i < j$. In the SBM, the value of connectivity probability $P_{i,j}$ depends on the communities the two endpoints i and j belong to. We assume $P_{i,j} = p$ if both nodes come from the same community and $P_{i,j} = q$ otherwise. There are k communities in the network. We denote $z \in [k]^n$, as the assignment vector, with z_i indicating the index of community the i -th node belongs to. Thus, the connectivity probability matrix P can be written as

$$P_{i,j} = B_{z_i, z_j}, \quad (2.1)$$

where $B \in [0, 1]^{k \times k}$ with diagonal entries as p and off-diagonal entries as q . That is, $B = q1_k1_k^T + (p - q)I_k$.

We consider a SBM with parameter space defined as follows,

$$\mathcal{Z}(n, k, \beta) \triangleq \left\{ z \in [k]^n : \min_{u \in [k]} |\{i : z_i = u\}| \geq \beta n/k - 1 \right\}, \quad (2.2)$$

where $\beta \geq 1$ to pose minimum community size. Here we allow β to be dependent on n . The “ -1 ” term is to avoid the extra constraint for the equal-community-size case when $\beta = 1$.

The goal is to estimate z from the observation A , with all the other parameters n, k, p, q, β known.

The existence and strength of community structure is determined by the difference between p and q . As we will show in later chapters, the signal-to-noise ratio essentially takes the form $(p - q)^2 / (p + q)$. To have a community structure, p can be either greater or smaller than q . Nevertheless, in this thesis we restrict to the $p \geq q$ case, i.e., the within-community probabilities are larger than the between-communities probabilities, as in reality individuals from the same community are often more likely to be connected. However almost identical results hold for the opposite case with $p < q$. Throughout this thesis, we assume $c_0/n < q \leq p \leq 1 - c_0$, where $0 < c_0 < 1$ can be any small constant, allowing the network to be from very sparse to very dense.

Inhomogeneous Stochastic Block Model. The aforementioned SBM might be restrictive, as the within-community and cross-community connection probabilities are homogeneous, in the sense that they are either p, q . A slightly more general case is inhomogeneous SBM, where we allow more flexibility in the B matrix.

In inhomogeneous SBM, we allow the diagonal entries of B (within-community connection probabilities) to be greater than p , and the off-diagonal entries of B (cross-community connection probabilities) to be smaller than q . Its formal definition is given as follows, we have $P_{i,j} = B_{z_i, z_j}, \forall i < j$ same as homogeneous SBM, but with a larger parameter space defined as

$$\Theta(n, k, \beta, p, q) = \left\{ (z, B) : \min_{u \in [k]} |\{i : z_i = u\}| \geq \beta n/k - 1, B_{u,u} \geq p, \forall u \in [k], \text{ and } B_{u,v} \leq q, \forall u \neq v \right\},$$

where we also assume $p > q$.

For simplicity, in this thesis, we study the community detection problem mainly under the regular SBM $\mathcal{Z}(n, k, \beta)$, with some extensions to the inhomogeneous SBM $\Theta(n, k, \beta, p, q)$.

Chapter 3

Minimax Rates

We will establish minimax rates for SBM under the parameter space $\mathcal{Z}(n, k, \beta)$ defined in Equation (2.2). But before that, we first introduce the concept of misclassification proportion, which will be used as the loss function, and a key quantity I , which is the signal-to-noise ratio.

Misclassification Proportion. There is an identifiability issue for any $z \in \mathcal{Z}(n, k, \beta)$, as the labels $1, 2, \dots, k$ are only identifiable up to a global shift. That is, z and $\rho \circ z$ give the same partition of the network, where $\rho : [k] \rightarrow [k]$ is any permutation over $[k]$. Consequently, for any $z, z' \in \mathcal{Z}(n, k, \beta)$, we define their distance as

$$\ell(z', z) = \frac{1}{n} \min_{\rho} \|\rho \circ z' - z\|_0 = \frac{1}{n} \min_{\rho} \sum_{i=1}^n \mathbb{I} \{ \rho(z'_i) = z_i \}, \quad (3.1)$$

where the minimization is over all the permutations on $[k]$. Note that $\|x - y\|_0$ for arbitrary vectors x, y is the same as their Hamming distance.

Signal-to-noise Ratio. We define a key quantity I as the Rényi divergence of order $1/2$ between two Bernoulli distributions $\text{Ber}(p)$ and $\text{Ber}(q)$, which has an explicit formula as

$$I = -2 \log \left(\sqrt{pq} + \sqrt{(1-p)(1-q)} \right). \quad (3.2)$$

To see that I can be interpreted as signal-to-noise ratio, by Proposition 4.1 we have $I = (1 + o(1))(\sqrt{p} - \sqrt{q})^2$ which is equal to $(p - q)^2/(p + q)$ up to a constant, when $p, q = o(1)$.

We present the minimax rates in Theorem 3.1 as follows.

Theorem 3.1. *Under the assumption $\beta n I / (k \log k) \rightarrow \infty$, there exists a sequence $\eta = o(1)$ that only depends on n, k, β, p, q , such that*

$$\min_{\hat{z}} \max_{z \in \mathcal{Z}(n, k, \beta)} \mathbb{E} \ell(\hat{z}, z) = \begin{cases} \exp(-(1 + \eta)nI/2), & k = 2; \\ \exp(-(1 + \eta)\beta n I / k), & k \geq 3. \end{cases} \quad (3.3)$$

Theorem 3.1 covers both dense and sparse networks. It holds for a wide range of possible values of p and q ranging from $1/n$ order to constant order. The number of communities k is allowed to grow fast. It can be as large as in the order of $n/\log n$ when the connectivity probability is nearly a constant order, in which each community contains an order of $\log n$ nodes. In addition, for finite number of communities, Theorem 3.1 shows $(p - q)^2/p \rightarrow \infty$ is a necessary and sufficient condition for consistent community detection, which implies consistency results in [37, 39]. It also recovers the strong consistency results in [26, 38], in which they additionally assume $p \asymp (\log n)/n$.

The fundamental hardness of estimating z is essentially the same as estimating label of one single node (say, z_1), assuming the labels of all the remaining nodes (i.e., z_2, z_3, \dots, z_n) are known. Investigating this *local* testing problem is crucially important, as it not only provides insights for the minimax rates, but late also inspires us algorithmically (see Chapter 4). After exploring this local problem in Section 3.1, we will obtain the lower bound of Theorem 3.1 by using a novel *global to local* scheme in Section 3.2. The minimax upper bound is obtained by maximum likelihood estimation (MLE), which will be provided in Section 3.3.

3.1 Hypothesis Testing for One Single Node

Let m_1, m_2 be arbitrary positive integers. Consider a network with $m_1 + m_2 + 1$ nodes and two underlying communities. Assume we know the labels of the last $m_1 + m_2$ nodes: $z_2, \dots, z_{m_1+1} = 1$ and $z_{m_1+2}, \dots, z_{m_1+m_2+1} = 2$. The task is to estimate z_1 which has two possibilities $z_1 = 1$ or $z_1 = 2$. The only difference in the adjacency matrix A is the first row. If $z_1 = 1$ then the first half of $\{A_{1,i}\}_{i=2}^{m_1+m_2+1}$ is $\text{Ber}(p)$ and the second half is $\text{Ber}(q)$; otherwise, the first half is $\text{Ber}(q)$ and the second half is $\text{Ber}(p)$.

The aforementioned testing problem can be formulated as follows. Let $\{X_i\}_{i=1}^{m_1}, \{Y_i\}_{i=1}^{m_2}$ be independent variables. We have the following two hypotheses.

$$\begin{aligned} H_0 &: \{X_i\}_{i=1}^{m_1} \stackrel{iid}{\sim} \text{Ber}(q), \{Y_i\}_{i=1}^{m_2} \stackrel{iid}{\sim} \text{Ber}(p); \\ H_1 &: \{X_i\}_{i=1}^{m_1} \stackrel{iid}{\sim} \text{Ber}(p), \{Y_i\}_{i=1}^{m_2} \stackrel{iid}{\sim} \text{Ber}(q). \end{aligned}$$

We have Lemma 3.1 to lower bound the summation of its Type I and Type II errors.

Lemma 3.1. *Let ϕ be any procedure based on $\{X_i\}_{i=1}^{m_1}, \{Y_i\}_{i=1}^{m_2}$. Define $m = \max\{m_1, m_2\}$. There exists a positive sequence $\eta = o(1)$ and a positive constant c such that*

$$\min_{\phi} \left(\frac{1}{2} \mathbb{P}_{H_0}(\phi = 1) + \frac{1}{2} \mathbb{P}_{H_1}(\phi = 0) \right) \geq \begin{cases} \exp(-(1 + \eta)mI), & \text{if } mI \rightarrow \infty; \\ c, & \text{if } mI = O(1). \end{cases}$$

Proof. Define another hypothesis testing problem as

$$\begin{aligned} \bar{H}_0 &: \{X_i\}_{i=1}^m \stackrel{iid}{\sim} \text{Ber}(q), \{Y_i\}_{i=1}^m \stackrel{iid}{\sim} \text{Ber}(p); \\ \bar{H}_1 &: \{X_i\}_{i=1}^m \stackrel{iid}{\sim} \text{Ber}(p), \{Y_i\}_{i=1}^m \stackrel{iid}{\sim} \text{Ber}(q), \end{aligned}$$

with two equal-length vectors. We have

$$\min_{\phi} \left(\frac{1}{2} \mathbb{P}_{H_0}(\phi = 1) + \frac{1}{2} \mathbb{P}_{H_1}(\phi = 0) \right) \geq \min_{\phi} \left(\frac{1}{2} \mathbb{P}_{\bar{H}_0}(\phi = 1) + \frac{1}{2} \mathbb{P}_{\bar{H}_1}(\phi = 0) \right)$$

Note that this is a Bayes risk with respect to a zero-one loss. Let $\hat{\phi}$ be the optimal procedure,

than $\hat{\phi}$ must be the mode of the posterior distribution. Since the prior is uniform, it immediately implies that $\hat{\phi}$ must be the likelihood ratio test. Note the likelihood function under \bar{H}_0 is

$$\begin{aligned} f_{\bar{H}_0} &= \prod_{i=1}^m q^{X_i} (1-q)^{1-X_i} \prod_{i=1}^m p^{Y_i} (1-p)^{1-Y_i} \\ &= \exp \left(\sum_{i=1}^m \left(X_i \log \frac{q}{1-q} + \log(1-q) \right) + \sum_{i=1}^m \left(Y_i \log \frac{p}{1-p} + \log(1-p) \right) \right), \end{aligned}$$

and similarly under \bar{H}_1 is

$$f_{\bar{H}_1} = \exp \left(\sum_{i=1}^m \left(X_i \log \frac{p}{1-p} + \log(1-p) \right) + \sum_{i=1}^m \left(Y_i \log \frac{q}{1-q} + \log(1-q) \right) \right).$$

As a consequence, $\hat{\phi}$ can be written explicitly as

$$\begin{aligned} \hat{\phi} &= \mathbb{I} \{ f_{\bar{H}_1} \geq f_{\bar{H}_0} \} \\ &= \mathbb{I} \left\{ \sum_{i=1}^m \left(X_i \log \frac{p(1-q)}{q(1-p)} - \log \frac{1-q}{1-p} \right) + \sum_{i=1}^m \left(Y_i \log \frac{q(1-p)}{p(1-q)} + \log \frac{1-q}{1-p} \right) \geq 0 \right\}. \end{aligned}$$

It can be further simplified into

$$\hat{\phi} = \mathbb{I} \left\{ \sum_{i=1}^m X_i - \sum_{i=1}^m Y_i \geq 0 \right\}. \quad (3.4)$$

Define $W = U - V$ where $U \sim \text{Ber}(q)$, $V \sim \text{Ber}(q)$ and U, V are independent. Let $\{W_i\}_{i=1}^m$ be i.i.d. copies of W . We have

$$\begin{aligned} \min_{\phi} \left(\frac{1}{2} \mathbb{P}_{H_0}(\phi = 1) + \frac{1}{2} \mathbb{P}_{H_1}(\phi = 0) \right) &\geq \frac{1}{2} \mathbb{P}_{\bar{H}_0} \left(\sum_{i=1}^m X_i - \sum_{i=1}^m Y_i \geq 0 \right) + \frac{1}{2} \mathbb{P}_{\bar{H}_1} \left(\sum_{i=1}^m X_i - \sum_{i=1}^m Y_i < 0 \right) \\ &= \mathbb{P} \left(\sum_{i=1}^m W_i \geq 0 \right). \end{aligned}$$

The proof is complete by using Lemma 3.2 to lower bound the RHS of above inequality. \blacksquare

The establishment of Lemma 3.2 mainly follows that of Cramer-Chernoff Theorem [48]. The general Cramer-Chernoff Theorem gives a lower bound for the tail probability that

the sum of random variables deviates from its mean. Usually it is for the case where these random variables are from a distribution independent of the sample size. In our setting we allow p and q to depend on n . We refer readers to [52] for its detailed proof.

Lemma 3.2 (Lemma 5.2 of [52]). *Assume $\{X_i\}_{i=1}^m \stackrel{iid}{\sim} \text{Ber}(q)$, $\{Y_i\}_{i=1}^m \stackrel{iid}{\sim} \text{Ber}(p)$ which are independent of each other. If $mI \rightarrow \infty$, there exists a sequence $\eta = o(1)$ such that*

$$\mathbb{P}\left(\sum_{i=1}^m (X_i - Y_i) \geq 0\right) \geq \exp(-(1 + \eta)mI).$$

In addition, if $mI = O(1)$, then $\mathbb{P}(\sum_{i=1}^m (X_i - Y_i) \geq 0) \geq c$ for some constant $c > 0$.

3.2 Minimax Lower Bound

The similarity of forms between Theorem 3.1 and Lemma 3.1 is hard not to notice. When taking $m = n/2$ or $\beta n/k$, the lower bound in Lemma 3.1 matches with the minimax rate in Theorem 3.1. Actually, as we will show in the section, the key to establish the minimax lower bound is to follow a novel *global to local* scheme: reducing the *global* community detection problem to a *local* hypothesis testing problem.

Dealing with $\ell(\hat{z}, z)$ directly is difficult and intimidating, as $\ell(\cdot, \cdot)$ involves with a minimization over all permutations, unless we are able to decoupling it into estimation errors of individual nodes. We use \mathcal{Z} short for $\mathcal{Z}(n, k, \beta)$ defined in Equation (2.2) for simplicity.

Case $k \geq 3$. Let $z^* \in \mathcal{Z}$ be an arbitrary assignment vector, such that

$$|\{i : z_i^* = 1\}| = |\{i : z_i^* = 2\}| = \beta n/k - 1.$$

We fix a set $T \in [n]$ such that $|T \cap \{i : z_i^* = u\}| = \delta \beta n/k^2, \forall u \in [k]$, where $\delta = o(1)$ is some sequence whose value will be specified later. We define a subspace of \mathcal{Z} as follow

$$\mathcal{Z}^* = \{z \in \mathcal{Z} : z_i = z_i^*, \forall i \in T^C\}.$$

In this way, for any $z, z' \in \mathcal{Z}^*$, they differ at most $|T| = \delta\beta n/k$ nodes, which implies that

$$\ell(z, z') = \frac{1}{n} \|z - z'\|_0 = \frac{1}{n} \sum_{i \in T} \mathbb{I}\{z_i = z'_i\}.$$

So we have

$$\inf_{\hat{z}} \sup_{z \in \mathcal{Z}} \mathbb{E}\ell(\hat{z}, z) \geq \inf_{\hat{z}} \sup_{z \in \mathcal{Z}^*} \mathbb{E}\ell(\hat{z}, z) = \inf_{\hat{z}} \sup_{z \in \mathcal{Z}^*} \frac{1}{n} \mathbb{E} \sum_{i \in T} \mathbb{I}\{\hat{z}_i = z_i\}.$$

Since minimax risk is lower bounded by Bayes risk, we have

$$\inf_{\hat{z}} \sup_{z \in \mathcal{Z}} \mathbb{E}\ell(\hat{z}, z) \geq \frac{1}{n|\mathcal{Z}^*|} \inf_{\hat{z}} \sum_{z \in \mathcal{Z}^*} \mathbb{E} \sum_{i \in T} \mathbb{I}\{\hat{z}_i = z_i\}.$$

Without loss of generality, suppose $1 \in T$. Due to symmetry, $\sum_{z \in \mathcal{Z}^*} \mathbb{E}\mathbb{I}\{\hat{z}_i = z_i\} = \sum_{z \in \mathcal{Z}^*} \mathbb{E}\mathbb{I}\{\hat{z}_1 = z_1\}$ for all $i \in T$. We have

$$\begin{aligned} \inf_{\hat{z}} \sup_{z \in \mathcal{Z}} \mathbb{E}\ell(\hat{z}, z) &\geq \frac{|T|}{n|\mathcal{Z}^*|} \inf_{\hat{z}} \sum_{z \in \mathcal{Z}^*} \mathbb{E}\mathbb{I}\{\hat{z}_1 = z_1\} \\ &\geq \frac{\delta\beta}{k|\mathcal{Z}^*|} \inf_{\hat{z}} \sum_{z \in \mathcal{Z}^*} \mathbb{P}_z(\hat{z}_1 = z_1). \end{aligned}$$

Here we have the subscript $\mathbb{P}_z(\hat{z}_1 = z_1)$ in \mathbb{P}_z to avoid confusion among different probability measures. We partition \mathcal{Z}^* into disjoint subsets $\mathcal{Z}^* = \cup_{u=1}^k \mathcal{Z}_u^*$ where $\mathcal{Z}_u^* = \{z \in \mathcal{Z}^* | z_1 = u\}$. Again due to symmetry, $|\mathcal{Z}_u^*|$ are equal to each other for all $k \in [k]$. Then

$$\inf_{\hat{z}} \sup_{z \in \mathcal{Z}} \mathbb{E}\ell(\hat{z}, z) \geq \frac{\delta\beta}{k|\mathcal{Z}^*|} \inf_{\hat{z}} \sum_{z \in \mathcal{Z}_1^* \cup \mathcal{Z}_2^*} \mathbb{P}_z(\hat{z}_1 = z_1) \geq \frac{\delta\beta}{k^2|\mathcal{Z}_1^*|} \inf_{\hat{z}} \sum_{z \in \mathcal{Z}_1^* \cup \mathcal{Z}_2^*} \mathbb{P}_z(\hat{z}_1 = z_1).$$

We are going to pair elements in \mathcal{Z}_1^* and \mathcal{Z}_2^* so that each pair only differs at the first node. For any $z \in \mathcal{Z}_1^*$ we denote z_{-1} to be remaining part of z after excluding the first node. We define a subspace of \mathcal{Z}_1^* as

$$\bar{\mathcal{Z}}_1^* = \{z \in \mathcal{Z}_1^* : (2, z_{-1}) \in \mathcal{Z}_2^*\}.$$

It turns out

$$\begin{aligned} \frac{|\bar{\mathcal{Z}}_1^*|}{|\mathcal{Z}_1^*|} &= 1 - \frac{\left| \left\{ z \in \mathcal{Z}_1^* \mid |\{i : z_i = 1\}| = \beta n/k - 1 \right\} \right|}{|\mathcal{Z}_1^*|} \\ &\geq 1 - \frac{\left| \left\{ z \in \mathcal{Z}_1^* \mid |\{i : z_i = 1\}| = \beta n/k - 1 \right\} \right|}{\left| \left\{ z \in \mathcal{Z}_1^* \mid |\{i : z_i = 1\}| = \beta n/k - 1 \right\} \right| + \left| \left\{ z \in \mathcal{Z}_1^* \mid |\{i : z_i = 1\}| = \beta n/k \right\} \right|}. \end{aligned}$$

Now we are going to build the connections between two cardinalities $|\bar{\mathcal{Z}}_{1,1}^*|$ and $|\bar{\mathcal{Z}}_{1,2}^*|$, where $\bar{\mathcal{Z}}_{1,1}^* = \left\{ z \in \mathcal{Z}_1^* \mid |\{i : z_i = 1\}| = \beta n/k - 1 \right\}$ and $\bar{\mathcal{Z}}_{1,2}^* = \left\{ z \in \mathcal{Z}_1^* \mid |\{i : z_i = 1\}| = \beta n/k \right\}$. For any $z \in \bar{\mathcal{Z}}_{1,2}^*$, for any $i \in \{j \in T : j \neq 1, z_j = 1\}$ (the cardinality of which is $\delta\beta n/k^2$), we can construct $z' \in |\bar{\mathcal{Z}}_{1,1}^*|$ by letting $z'_j = z_j, \forall j \neq i$ and $z'_i \neq z_i$. Hence there are $(\delta\beta n/k^2)(k-1)$ different $z' \in |\bar{\mathcal{Z}}_{1,1}^*|$ that can be generated by one single $z \in |\bar{\mathcal{Z}}_{1,2}^*|$. One the other hand, for any $z' \in |\bar{\mathcal{Z}}_{1,1}^*|$, it can be generated by the aforementioned procedure by at least $(\delta\beta n/k - \delta\beta n/k^2)$ (recall that $|T| = \delta\beta n/k$) different $z \in |\bar{\mathcal{Z}}_{1,2}^*|$. This leads to

$$\frac{|\bar{\mathcal{Z}}_{1,1}^*|}{|\bar{\mathcal{Z}}_{1,2}^*|} \leq \frac{(\delta\beta n/k^2)(k-1)}{(\delta\beta n/k - \delta\beta n/k^2)} = 1.$$

Thus,

$$\frac{|\bar{\mathcal{Z}}_1^*|}{|\mathcal{Z}_1^*|} \geq 1 - \frac{|\bar{\mathcal{Z}}_{1,1}^*|}{|\bar{\mathcal{Z}}_{1,1}^*| + |\bar{\mathcal{Z}}_{1,2}^*|} \geq \frac{1}{2}.$$

We have

$$\begin{aligned} \inf_{\hat{z}} \sup_{z \in \mathcal{Z}} \mathbb{E} \ell(\hat{z}, z) &\geq \frac{\delta\beta}{2k^2 |\bar{\mathcal{Z}}_1^*|} \inf_{\hat{z}} \sum_{z \in \mathcal{Z}_1^* \cup \mathcal{Z}_2^*} \mathbb{P}_z(\hat{z}_1 = z_1) \\ &\geq \frac{\delta\beta}{2k^2 |\bar{\mathcal{Z}}_1^*|} \inf_{\hat{z}} \sum_{z \in \bar{\mathcal{Z}}_1^*} (\mathbb{P}_z(\hat{z}_1 = z_1) + \mathbb{P}_{(2, z_{-1})}(\hat{z}_1 = z_1)) \\ &\geq \frac{\delta\beta}{k^2 |\bar{\mathcal{Z}}_1^*|} \sum_{z \in \bar{\mathcal{Z}}_1^*} \inf_{\hat{z}} \left(\frac{1}{2} \mathbb{P}_z(\hat{z}_1 = z_1) + \frac{1}{2} \mathbb{P}_{(2, z_{-1})}(\hat{z}_1 = z_1) \right) \\ &= \frac{\delta\beta}{k^2 |\bar{\mathcal{Z}}_1^*|} \sum_{z \in \bar{\mathcal{Z}}_1^*} \inf_{\hat{z}_1} \left(\frac{1}{2} \mathbb{P}_z(\hat{z}_1 = z_1) + \frac{1}{2} \mathbb{P}_{(2, z_{-1})}(\hat{z}_1 = z_1) \right). \end{aligned}$$

We can use the arguments in Section 3.1 and Lemma 3.1 for lower bounding Bayes risk for two-point hypothesis testing problem. For any $z \in \bar{\mathcal{Z}}_1^*$, we have $\max_{u=1,2} \{|\{i \in [n] : z_i =$

$u\}}\} \leq \beta n/k + \delta \beta n/k$. Thus

$$\frac{1}{2} \mathbb{P}_z(\hat{z}_1 = z_1) + \frac{1}{2} \mathbb{P}_{(2, z_{-1})}(\hat{z}_1 = z_1) \geq \exp(-(1 + \eta')(1 + \delta)\beta n I/k),$$

for some positive sequence $\eta' = o(1)$, if $\beta n I/k \rightarrow \infty$. Otherwise it is lower bounded by some constant $c > 0$. Then we have

$$\inf_{\hat{z}} \sup_{z \in \mathcal{Z}} \mathbb{E} \ell(\hat{z}, z) \geq \frac{\delta \beta}{k^2} \exp\left(- (1 + \eta')(1 + \delta) \frac{\beta n I}{k}\right).$$

Under the assumption $\beta n I/(k \log k) \rightarrow \infty$ and $\log(1/\delta) = o(\beta n I/k)$, there exists a positive sequence $\eta = o(1)$ such that

$$\inf_{\hat{z}} \sup_{z \in \mathcal{Z}} \mathbb{E} \ell(\hat{z}, z) \geq \exp(-(1 + \eta)\beta n I/k).$$

Case $k = 2$. The least favorable situation is when two communities are almost equal sized. We construct \mathcal{Z}^* analogous as the case $k \geq 3$, but instead require the community sizes to be almost equal. The proof is almost identical to the case $k \geq 3$, and hence is omitted in this thesis.

3.3 Minimax Upper Bound

The minimax upper bound can be achieved by maximum likelihood estimator (MLE). For any $z \in \mathcal{Z}(n, k, \beta)$, its likelihood function takes a form as

$$\begin{aligned} f(z; A) &= \prod_{i < j, z_i = z_j} p^{A_{i,j}} (1-p)^{1-A_{i,j}} \prod_{i < j, z_i \neq z_j}^m q^{A_{i,j}} (1-q)^{1-A_{i,j}} \\ &= \exp\left(\sum_{i < j, z_i = z_j} \left(A_{i,j} \log \frac{p}{1-p} + \log(1-p) \right) + \sum_{i < j, z_i \neq z_j} \left(A_{i,j} \log \frac{q}{1-q} + \log(1-q) \right) \right) \\ &= C \exp\left(\sum_{i < j, z_i = z_j} \left(A_{i,j} \log \frac{p(1-q)}{q(1-p)} - \log \frac{1-q}{1-p} \right) \right) \end{aligned}$$

for some C independent of z . Here the third equation is due to the fact $\prod_{i < j, z_i = z_j} A_{i,j} + \prod_{i < j, z_i \neq z_j} A_{i,j} = \sum_{i < j} A_{i,j}$ does not depend on the choice of z . Thus the MLE has an expression as

$$\begin{aligned}
\hat{z}^{\text{MLE}} &= \operatorname{argmax}_{z \in \mathcal{Z}(n, k, \beta)} \log f(z; A) \\
&= \operatorname{argmax}_{z \in \mathcal{Z}(n, k, \beta)} \sum_{i < j, z_i = z_j} \left(A_{i,j} \log \frac{p(1-q)}{q(1-p)} - \log \frac{1-q}{1-p} \right) \\
&= \operatorname{argmax}_{z \in \mathcal{Z}(n, k, \beta)} \sum_{i < j, z_i = z_j} (A_{i,j} - \lambda),
\end{aligned} \tag{3.5}$$

where λ is defined as

$$t = \frac{1}{2} \log \frac{p(1-q)}{q(1-p)}, \quad \lambda = \frac{1}{2t} \log \frac{1-q}{1-p}. \tag{3.6}$$

Now we are going to prove \hat{z}^{MLE} is a minimax estimator. Denote $z^* \in \mathcal{Z}(n, k, \beta)$ be the ground truth where A is generated from. We will show it is unlikely for \hat{z}^{MLE} to be far away from z^* , which is equivalent to show it is unlikely some z not in a neighborhood of z^* has likelihood $f(z; A)$ greater than $f(z^*; A)$. The main technique we use to prove it is union bound with chaining.

Note that

$$\begin{aligned}
\mathbb{I}\{f(z; A) > f(z^*; A)\} &\iff \mathbb{I}\left\{ \sum_{i < j, z_i = z_j} (A_{i,j} - \lambda) > \sum_{i < j, z_i^* = z_j^*} (A_{i,j} - \lambda) \right\} \\
&\iff \mathbb{I}\left\{ \sum_{i < j, z_i = z_j, z_i^* \neq z_j^*} A_{i,j} - \sum_{i < j, z_i \neq z_j, z_i^* = z_j^*} A_{i,j} > \lambda(\gamma(z; z^*) - \alpha(z; z^*)) \right\},
\end{aligned}$$

where

$$\alpha(z; z^*) = |\{(i, j) : i < j, z_j \neq z_j^*, z_i^* = z_j^*\}|, \tag{3.7}$$

$$\text{and } \gamma(z; z^*) = |\{(i, j) : i < j, z_j = z_j^*, z_i^* \neq z_j^*\}|. \tag{3.8}$$

We have Proposition 3.1 to upper bound the probability $\mathbb{P}(f(z; A) \geq f(z^*; A))$, which is an

immediate consequence of Chernoff bound.

Proposition 3.1. *Let $z, z^* \in \mathcal{Z}(n, k, \beta)$ with z^* being the ground truth. Then*

$$\mathbb{P}(f(z; A) \geq f(z^*; A)) \leq \exp(-(\alpha(z; z^*) + \gamma(z; z^*))I/2).$$

The following Proposition provides control on $\alpha(z; z^*) + \gamma(z; z^*)$.

Proposition 3.2. *For any $z, z^* \in \mathcal{Z}(n, k, \beta)$. Denote $m = n\ell(z, z^*)$, we have*

$$\alpha(z; z^*) + \gamma(z; z^*) \geq \begin{cases} m(n - m), & k = 2; \\ \beta nm / (16k), & k \geq 3. \end{cases}$$

In addition, when $k \geq 3$ and $m \leq \beta n / (2k)$, we have $\alpha(z; z^) + \gamma(z; z^*) \geq 2(\beta nm / k - m^2)$.*

Proposition 3.3 provides a control of cardinality as we will later implement union bound. It is worthwhile pointing out that, we should not use the cardinality of $\{z \in \mathcal{Z}(n, k, \beta) : n\ell(z, z^*) = m\}$, which is too large due to counting assignments equivalent under permutation.

Proposition 3.3. *The cardinality of equivalent class that has distance m from z^* is upper bounded as follows,*

$$\left| \left\{ \Gamma : \exists z \in \Gamma \text{ s.t. } n\ell(z, z^*) = m \right\} \right| \leq \min \left\{ \left(\frac{enk}{m} \right)^m, k^n \right\},$$

where $0 < m < n$ is a positive integer.

Define P_m as

$$P_m = \mathbb{P}(\exists z \in \mathcal{Z}(n, k, \beta) : n\ell(z, z^*) = m \text{ and } f(z; A) \geq f(z^*; A)). \quad (3.9)$$

By union bound we have

$$\mathbb{P}_m \leq \left| \left\{ \Gamma : \exists z \in \Gamma \text{ s.t. } n\ell(z, z^*) = m \right\} \right| \max_{z: n\ell(z, z^*)=m} \mathbb{P}(f(z; A) \geq f(z^*; A)).$$

In the following we consider the case $k \geq 3$. The proof for the case when $k = 2$ is similar and hence omitted. Note that we have $\frac{\beta n I}{k \log k} \rightarrow \infty$. We consider three scenarios as follows.

1) If $\beta n I/k > (1 + \epsilon) \log n$, define $m_0 = 1$ and $m' = \epsilon \beta n/k$. Then $P_1 \leq n \exp(-(\beta n/k - 1)I)$. Denote $R = n \exp(-(\beta n/k - 1)I)$. We have

$$P_m \leq \begin{cases} \left(\frac{\epsilon n k}{2}\right)^m \exp(-(\beta n m/k - m^2)I) \leq R n^{-\epsilon(m-1)/6}, & \text{for } m_0 < m \leq m' \\ \left(\frac{\epsilon n k}{m'}\right)^m \exp(-\frac{\beta n m I}{32k}) \leq R \exp(-\frac{\beta n(m-4)I}{64k}), & \text{for } m' < m \leq n/2. \end{cases}$$

Then $n \mathbb{E} \ell(\hat{z}, z^*) \leq \sum_{m=1}^{n/2} m P_m = (1 + o(1))R$.

2) If $\beta n I/k < (1 - \epsilon) \log n$, define $m_0 = n \exp(-(1 - e^{-\epsilon n I/2})\beta n I/k)$ and $m' = n \exp(-\beta n I/8k)$.

We have

$$P_m \leq \begin{cases} \left(\frac{\epsilon n k}{m_0}\right)^m \exp(-(\beta n m/k - m^2)I) = \exp(-e^{-\frac{\epsilon n I}{2}} \frac{\beta n m I}{4k}), & \text{for } m_0 < m \leq m' \\ \left(\frac{\epsilon n k}{m'}\right)^m \exp(-\frac{\beta n m I}{32k}) \leq \exp(-\frac{\beta n m I}{64k}), & \text{for } m' < m \leq n/2. \end{cases}$$

Then $\mathbb{E} \ell(\hat{z}, z^*) \leq m_0/n + \sum_{m > m_0}^{n/2} P_m = (1 + o(1))m_0/n$.

3) If $\frac{\beta n I}{k \log n} \rightarrow 1$, there exists a positive sequence $w \rightarrow 0$ such that $|\frac{\beta n I}{k \log n} - 1| \ll w$ and $\frac{1}{\sqrt{\log n}} \leq w$. Define $m_0 = n \exp(-(1 - w)\beta n I/k)$ and $m' = w^2 n$.

$$P_m \leq \begin{cases} \left(\frac{\epsilon n k}{m_0}\right)^m \exp(-(\beta n m/k - m^2)I) \leq \exp(-\frac{w \beta n m I}{2k}), & \text{for } m_0 < m \leq m' \\ \left(\frac{\epsilon n k}{m'}\right)^m \exp(-\frac{\beta n m I}{32k}) \leq \exp(-\frac{\beta n m I}{4k}), & \text{for } m' < m \leq n/2. \end{cases}$$

Then $\mathbb{E} \ell(\hat{z}, z^*) \leq m_0/n + \sum_{m > m_0}^{n/2} P_m = (1 + o(1))m_0/n$.

3.4 Extension

Theorem 3.1 can be generalized to inhomogeneous SBM $\Theta(n, k, \beta, p, q)$ with the same min-max rate holds.

Theorem 3.2. *Under the assumption $\beta n I/(k \log k) \rightarrow \infty$, there exists a sequence $\eta = o(1)$*

that only depends on n, k, β, p, q , such that

$$\min_{\hat{z}} \max_{(z, B) \in \Theta(n, k, \beta, p, q)} \mathbb{E} \ell(\hat{z}, z) = \begin{cases} \exp(-(1 + \eta)nI/2), & k = 2; \\ \exp(-(1 + \eta)\beta nI/k), & k \geq 3. \end{cases}$$

Proof. The proof of lower bound is trivial, noticing that $\mathcal{Z}(n, k, \beta)$ serves as the least favorable case for the inhomogeneous SBM. For the upper bound, we use the same procedure defined as in Equation (3.5). This is no longer the MLE for $\Theta(n, k, \beta, p, q)$, but still minimax optimal. Its proof is identical. \blacksquare

3.5 Proof of Propositions

Proof of Proposition 3.1. We use α, γ instead of $\alpha(z; z^*), \gamma(z; z^*)$ for simplicity. Note that

$$\mathbb{P}(f(z; A) \geq f(z^*; A)) = \mathbb{P}\left(\sum_{i=1}^{\gamma} X_i - \sum_{i=1}^{\alpha} Y_i \geq \lambda(\gamma - \alpha)\right),$$

where $\{X_i\}_{i=1}^{\gamma} \stackrel{iid}{\sim} \text{Ber}(q)$ and $\{Y_i\}_{i=1}^{\alpha} \stackrel{iid}{\sim} \text{Ber}(p)$ and they are independent of each other. By Chernoff bound, we have

$$\mathbb{P}(f(z; A) \geq f(z^*; A)) \leq \exp(-t\lambda(\gamma - \alpha)) (\mathbb{E} \exp(tX))^{\gamma} (\mathbb{E} \exp(-tY))^{\alpha},$$

where t is defined as in Equation (3.6). Simple algebra with Proposition 3.4 immediately leads to the desired result. \blacksquare

Proposition 3.4. *Assume $0 < q < p < 1$. Let $X \sim \text{Ber}(q)$ and $Y \sim \text{Ber}(p)$. Recall the definition $\lambda = \log \frac{1-q}{1-p} / \log \frac{p(1-q)}{q(1-p)}$, $t = \frac{1}{2} \log \frac{p(1-q)}{q(1-p)}$ and $I = -2 \log[\sqrt{pq} + \sqrt{(1-p)(1-q)}]$. Then the following two equations hold*

$$e^{t\lambda} = \left(\frac{\mathbb{E} e^{tX}}{\mathbb{E} e^{-tY}}\right)^{\frac{1}{2}}, \text{ and } \mathbb{E} e^{tX} \mathbb{E} e^{-tY} = \exp(-I). \quad (3.10)$$

Proof. The proof is straightforward and all by calculation. Note that $\mathbb{E} \exp(tX) = qe^t + 1 - q$

and $\mathbb{E} \exp(tY) = pe^t + 1 - p$. We can easily obtain

$$\mathbb{E} e^{tX} \mathbb{E} e^{-tY} = (qe^t + 1 - q)(pe^{-t} + 1 - p) = (\sqrt{pq} + \sqrt{(1-p)(1-q)})^2 = \exp(-I).$$

We can justify the first part of Equation (3.10) in a similar way. ■

Proof of Proposition 3.2. The case when $k = 2$ is straightforward, noting that m has to be smaller than $n/2$ by the definition of $\ell(\cdot, \cdot)$. We mainly focus on $k \geq 3$ case.

We use α, γ instead of $\alpha(z; z^*), \gamma(z; z^*)$ for simplicity. Without loss of generality we assume $\|z - z^*\|_0 = n\ell(z, z^*)$. Define $\mathcal{C}_u = \{i : z_i^* = u\}$ and $L_{u,v} = \sum_{i \in \mathcal{C}_u} \mathbb{I}\{z_i = v\}$. We have the equality $\sum_v L_{u,v} = |\mathcal{C}_u|$ and also

$$\alpha = \frac{1}{2} \sum_u \left[|\mathcal{C}_u|^2 - \sum_w L_{u,w}^2 \right] = \frac{1}{2} \sum_u \sum_{w \neq u} L_{u,w} L_{u,w'}$$

and $\gamma = \frac{1}{2} \sum_{u \neq v} \sum_w L_{u,w} L_{v,w}$.

We define $[k]$ into two disjoint subsets S_1 and S_2 where

$$S_1 = \left\{ u \in [k] : \forall v \neq u, L_{u,v} \leq \frac{3}{4} |\mathcal{C}_u| \right\},$$

and $S_2 = \left\{ i \in [k] : \exists v \neq u, L_{u,v} > \frac{3}{4} |\mathcal{C}_u| \right\}$.

Define $L_u = \sum_{v \neq u} L_{u,v}$. For any $u \in S_1$, if $L_{u,u} \geq |\mathcal{C}_u|/4$, we have $|\mathcal{C}_u|^2 - \sum_w L_{u,w}^2 \geq L_{u,u} L_u \geq |\mathcal{C}_u| L_u / 4$. If $L_{u,u} < \frac{1}{4} |\mathcal{C}_u|$ we have $|\mathcal{C}_u|^2 - \sum_w L_{u,w}^2 \geq \frac{3}{8} |\mathcal{C}_u|^2 \geq |\mathcal{C}_u| L_u / 4$ as well.

This leads to

$$\alpha \geq \frac{1}{2} \sum_{u \in S_1} \left[|\mathcal{C}_u|^2 - \sum_w L_{u,w}^2 \right] \geq \frac{1}{8} \sum_{u \in S_1} |\mathcal{C}_u| L_u.$$

For any $u \in S_2$ there exists a $v \neq u$ such that $L_{u,v} > \frac{3}{4} |\mathcal{C}_u|$. We must have $L_{u,u} + L_{v,v} \geq L_{u,v} + L_{v,u}$ otherwise $\|z - z^*\|_0 = n\ell(z, z^*)$ does not hold since we can switch the u -th and v -th columns of z to make $\|z - z^*\|_0$ smaller. Consequently, we have $L_{v,v} \geq L_u / 2$. So we

have $\sum_{u' \neq u} \sum_w L_{u,w} L_{u',w} \geq L_{u,v} L_{v,v} \geq 3|\mathcal{C}_u|L_u/8$. Then we have

$$\gamma \geq \frac{1}{2} \sum_{u \in \mathcal{S}_2} \sum_{u' \neq u} \sum_w L_{u,w} L_{u',w} \geq \frac{3}{8} \sum_{u \in \mathcal{S}_2} |\mathcal{C}_u|L_u.$$

Thus,

$$\alpha + \gamma \geq \frac{1}{16} \sum_u |\mathcal{C}_u|L_u \geq \frac{\beta n}{16k} \sum_u L_u \geq \frac{\beta n}{16k} \|z - z^*\|_1 = \frac{\beta n m}{16k}.$$

We have stronger result when $m \leq \beta n/(2k)$. Without loss of generality we assume $m = n\ell(z, z^*)$, which is equivalent to be stated as $m = \sum_i \mathbb{I}\{z_i \neq z_i^*\}$. Define $m_u = \sum_i \mathbb{I}\{z_i \neq z_i^*, z_i^* = u\}, \forall u \in [k]$ and $\gamma_u = |\{(i, j) : i < j, z_j = z_j^*, z_i^* \neq z_j^*\}|, \forall u \in [k]$. Thus

$$\gamma_u \geq |\{i : z_i = z_i^* = u\}| |\{i : z_i = u, z_i^* \neq u\}| \geq m_u(\beta n/k - m_u), \forall u \in [k].$$

Note that $\sum_u m_u = m$ and $\sum_u \gamma_u = \gamma$. We have

$$\alpha + \gamma \geq \sum_u m_u(\beta n/k - m_u) \geq \beta n m/k - m^2.$$

■

Proof of Proposition 3.3. Without loss of generality we assume that $\|z - z^*\|_0 = m$. Then z assigns m nodes with different values from z^* , and there are k possible values for each node. Thus

$$\left| \left\{ \Gamma : \exists z \in \Gamma \text{ s.t. } \|z - z^*\|_0 = m \right\} \right| \leq \binom{n}{m} k^m \leq \left(\frac{enk}{m} \right)^m.$$

In addition, since each node has at most k possible choices, we have a naive bound for the cardinality of Γ as $|\{\Gamma\}| \leq k^n$. ■

Chapter 4

Methodology

Starting with the proposal of a series of methodologies [21, 27, 32, 42], we have seen a large literature devoted to algorithmic solutions to uncovering community structure. There are two desired properties of any algorithm: 1) *computability*: whether it is a polynomial-time algorithm or not; 2) *optimality*: whether it has provable theoretical guarantee that matches with minimax rate. Unfortunately, the existing methods do not satisfy both properties simultaneously, and this motivates us to propose a computationally feasible algorithm for community detection in the SBM with adaptive minimax optimal performance.

4.1 Hypothesis Testing for One Single Node: Revisit

Now let us revisit the two-point hypothesis testing problem introduced in Section 3.1.

$$\begin{aligned} H_0 &: \{X_i\}_{i=1}^{m_1} \stackrel{iid}{\sim} \text{Ber}(q), \{Y_i\}_{i=1}^{m_2} \stackrel{iid}{\sim} \text{Ber}(p); \\ H_1 &: \{X_i\}_{i=1}^{m_1} \stackrel{iid}{\sim} \text{Ber}(p), \{Y_i\}_{i=1}^{m_2} \stackrel{iid}{\sim} \text{Ber}(q). \end{aligned}$$

By Neyman-Pearson lemma, the likelihood ratio test is optimal. Same as in Section 3.1, the two likelihoods take forms as

$$f_{H_0} = \exp \left(\sum_{i=1}^{m_1} \left(X_i \log \frac{q}{1-q} + \log(1-q) \right) + \sum_{i=1}^{m_2} \left(Y_i \log \frac{p}{1-p} + \log(1-p) \right) \right),$$

and similarly under \bar{H}_1 is

$$f_{H_1} = \exp \left(\sum_{i=1}^{m_1} \left(X_i \log \frac{p}{1-p} + \log(1-p) \right) + \sum_{i=1}^{m_2} \left(Y_i \log \frac{q}{1-q} + \log(1-q) \right) \right).$$

The $\hat{\phi}$ the likelihood ratio test takes a form as

$$\begin{aligned} \hat{\phi} &= \mathbb{I}\{f_{H_1} \geq f_{H_0}\} \\ &= \mathbb{I} \left\{ \sum_{i=1}^{m_1} \left(X_i \log \frac{p(1-q)}{q(1-p)} - \log \frac{1-q}{1-p} \right) \geq \sum_{i=1}^{m_2} \left(Y_i \log \frac{p(1-q)}{q(1-p)} - \log \frac{1-q}{1-p} \right) \right\} \\ &= \mathbb{I} \left\{ \sum_{i=1}^{m_1} X_i - \sum_{i=1}^{m_2} Y_i \geq \lambda(m_1 - m_2) \right\}, \end{aligned}$$

where λ is defined in Equation (3.6). Chernoff bound leads to the following Lemma 4.1, which matches with the lower bound (i.e., Lemma 3.1) when $m_1 = m_2$.

Lemma 4.1. *The likelihood ratio test $\hat{\phi}$ satisfies*

$$\frac{1}{2} \mathbb{P}_{H_0}(\phi = 1) + \frac{1}{2} \mathbb{P}_{H_1}(\phi = 0) \leq \exp(-(m_1 + m_2)I/2).$$

Proof. Let $\{X_i\}_{i=1}^{m_1} \stackrel{iid}{\sim} X$ and $\{Y_i\}_{i=1}^{m_2} \stackrel{iid}{\sim} Y$ and they are independent of each other, where $X \sim \text{Ber}(q)$ and $Y \sim \text{Ber}(p)$. Then

$$\mathbb{P}_{H_0}(\phi = 1) = \mathbb{P} \left(\sum_{i=1}^{m_1} X_i - \sum_{i=1}^{m_2} Y_i \geq \lambda(m_1 - m_2) \right),$$

By Chernoff bound, we have

$$\mathbb{P}_{H_0}(\phi = 1) \leq \exp(-t\lambda(m_1 - m_2)) (\mathbb{E} \exp(tX))^{m_1} (\mathbb{E} \exp(-tY))^{m_2},$$

where t is defined as in Equation (3.6). Simple algebra with Proposition 3.4 immediately leads to $\mathbb{P}_{H_0}(\phi = 1) \leq \exp(-(m_1 + m_2)I/2)$. We have same result for $\mathbb{P}_{H_1}(\phi = 0)$. ■

4.2 Spectral Clustering

Spectral clustering has been one of the most popular methods for community detection. It has been investigated by [12, 13, 17, 30, 31, 34, 35, 44, 46, 47, 50] with provable upper bound established. The spectral clustering usually takes two steps:

1. Eigendecomposition on A . Let u_1, u_2, \dots, u_n be the eigenvectors corresponding to eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. Denote $U = [u_1, u_2, \dots, u_k]$ be the eigenspace with k -leading eigenvectors.
2. Perform k -means on the rows of U (i.e., $\{U_{i,\cdot}\}_{i=1}^n$) to partition the n nodes.

Despite the fact that spectral clustering being easy to implement, its existing theoretical results (e.g., [35]) require strong assumptions on the parameters. The main techniques used are Davis-Khan Theorem and a sharp upper bound $\|A - P\|_{\text{op}}$. The former one needs a lower bound on the eigengap of P : this is trivial when $z^* \in \mathcal{Z}(n, k, \beta)$ but may not hold for the inhomogeneous SBM as B may be singular. The latter one needs the network to be dense, in the sense that $p, q = O((\log n)/n)$ otherwise the desired upper bound on $\|A - P\|_{\text{op}}$ no longer holds.

We propose a novel *low-rank based* spectral clustering in Algorithm 1. The additional truncation step (i.e., Step 1) makes it possible to have similar control on the operator norm even for the sparse network. In Step 2 we use a low rank approximation instead of the eigenspace to avoid the use of Davis-Khan Theorem, thus we have no requirement on the eigengap any more. The provable result for Algorithm 1 is given in Theorem 4.1.

Algorithm 1: Low-rank Based Spectral Clustering

Input: Adjacency matrix $A \in \{0, 1\}^{n \times n}$, number of communities k , parameter p

Output: Partition of the network z

- 1 Define $T(A) \in \{0, 1\}^{n \times n}$ by replacing the i th row and column of A whose row sum or column sum is larger than $20np$ by zeroes for each $i \in [n]$;
 - 2 Singular value decomposition (SVD) on $T(A)$ to obtain $T(A) = \sum_{i=1}^n \lambda_i u_i u_i^T$ with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. Let $\hat{P} = \sum_{i=1}^k \lambda_i u_i u_i^T$ be the rank- k approximation of $T(A)$;
 - 3 Perform k -means on the rows of \hat{P} (i.e., $\{\hat{P}_{i,\cdot}\}_{i=1}^n$) to obtain a partition of the network.
-

Theorem 4.1. Assume $z^* \in \mathcal{Z}(n, k, \beta)$. Let \hat{z} be the result from Algorithm 1. Under the assumption $\beta^2 n(p - q)^2 / (k^3 p) \rightarrow \infty$, with probability at least $1 - n^{-2}$, we have

$$\ell(\hat{z}, z^*) \leq \frac{ck^2 p}{\beta n(p - q)^2},$$

for some positive constant c .

Proof. The proof contains two parts: one for Step 1 and 2, and the other for the k -means in Step 3.

1) By the definition of \hat{P} , we have

$$\hat{P} = \operatorname{argmin}_{\operatorname{rank}(X) \leq k} \|T(A) - X\|_{\mathbb{F}}.$$

Define $P' = P + pI_n$ such that P' is a rank- k matrix and differs from P only by the diagonal entries. Thus, we have $\|T(A) - \hat{P}\|_{\mathbb{F}} \leq \|T(A) - P'\|_{\mathbb{F}}$. After rearrangement, we have

$$\begin{aligned} \|\hat{P} - P\|_{\mathbb{F}}^2 &\leq 2 \left| \langle \hat{P} - P', T(A) - P \rangle \right| + \|P' - P\|_{\mathbb{F}}^2 \\ &\leq 2 \|\hat{P} - P'\|_{\mathbb{F}} \sup_{\{X: \|X\|_{\mathbb{F}}=1, \operatorname{rank}(X) \leq 2k\}} |\langle X, T(A) - P \rangle| + \|P' - P\|_{\mathbb{F}}^2 \\ &\leq \frac{1}{4} \|\hat{P} - P'\|_{\mathbb{F}}^2 + 4 \sup_{\{X: \|X\|_{\mathbb{F}}=1, \operatorname{rank}(X) \leq 2k\}} |\langle X, T(A) - P \rangle|^2 + \|P' - P\|_{\mathbb{F}}^2 \\ &\leq \frac{1}{2} \|\hat{P} - P'\|_{\mathbb{F}}^2 + \frac{3}{2} \|P' - P\|_{\mathbb{F}}^2 + 4 \sup_{\{X: \|X\|_{\mathbb{F}}=1, \operatorname{rank}(X) \leq 2k\}} |\langle X, T(A) - P \rangle|^2. \end{aligned}$$

Therefore,

$$\|\hat{P} - P\|_{\mathbb{F}}^2 \leq 3 \|P' - P\|_{\mathbb{F}}^2 + 8 \sup_{\{X: \|X\|_{\mathbb{F}}=1, \operatorname{rank}(X) \leq 2k\}} |\langle X, T(A) - P \rangle|^2. \quad (4.1)$$

Apply singular value decomposition to X and we get $X = \sum_{l=1}^{2k} \sigma_l v_l v_l^T$. Then,

$$|\langle X, T(A) - P \rangle| \leq \sum_{l=1}^{2k} |\sigma_l| |v_l^T (T(A) - P) v_l| \leq \|T(A) - P\|_{\operatorname{op}} \sum_{l=1}^{2k} |\sigma_l| \leq \sqrt{2k} \|T(A) - P\|_{\operatorname{op}}.$$

By Lemma 4.2, we have $\|T(A) - P\|_{\operatorname{op}} \leq c\sqrt{np}$ with probability at least $1 - n^{-2}$ for some

constant $c > 0$. Hence, $\sup_{\{X: \|X\|_F=1, \text{rank}(X) \leq 2k\}} |\langle X, T(A) - P \rangle|^2 \leq c^2 np$, with probability at least $1 - n^{-2}$. Moreover, $\|P' - P\|_F^2 = np$. Using Equation (4.1), with probability at least $1 - n^{-2}$ we have

$$\|\hat{P} - P\|_F^2 \leq (3 + 16c^2k)np,$$

and consequently by triangle inequality $\|\hat{P} - P'\|_F \leq \|\hat{P} - P\|_F + \|P - P'\|_F$ we have

$$\|\hat{P} - P'\|_F^2 \leq c' knp,$$

for some constant c' .

2) By the definition of k -means, its output $\hat{z}, \{v_u\}_{u=1}^k$ satisfies

$$(\hat{z}, \{v_u\}_{u=1}^k) = \underset{z}{\operatorname{argmin}} \underset{\{v_u\}_{u=1}^k}{\operatorname{argmin}} \sum_{i=1}^n \left\| \hat{P}_{i,\cdot} - v_{z_i} \right\|^2.$$

Let $V \in \mathbb{R}^{n \times n}$ such that $V_{i,\cdot} = v_{z_i}, \forall i \in [n]$. Then we have

$$\|V - \hat{P}\|_F^2 \leq \|P' - \hat{P}\|_F^2 \leq c' knp.$$

Note that P' only have k unique rows which are well separated. If $z_i = z_j$ then $\|P'_{i,\cdot} - P'_{j,\cdot}\|^2 = 0$; and

$$\|P'_{i,\cdot} - P'_{j,\cdot}\|^2 \geq 2(p - q)^2 \beta n / k, \text{ for all } (i, j) \text{ such that } z_i \neq z_j. \quad (4.2)$$

Define

$$S = \left\{ i : \left\| V_{i,\cdot} - \hat{P}'_{j,\cdot} \right\|^2 \geq (p - q)^2 \beta n / (4k) \right\}.$$

Then

$$|S|(p - q)^2 \beta n / (2k) \leq \|V - \hat{P}\|_F^2,$$

which implies

$$n^{-1}|S| \leq \frac{4k^2 c' p}{\beta n (p - q)^2}.$$

We are going to show under the assumption $\beta^2 n (p - q)^2 / (k^3 p) \rightarrow \infty$ (which implies $|S| = o(\beta n / k)$), all the nodes in S^C will be correctly clustered. Define

$$C_u = \{i \in [n] : z_i^* = u, i \in S^C\}, \forall u \in [k].$$

We have the following arguments:

- For each $u \in [k]$, C_u cannot be empty, as $|C_u| \geq |\{i : z_i^* = u\}| - |S| > 0$.
- For each pair $u \neq v$, there cannot exist some $i \in C_u, j \in C_v$ such that $\hat{z}_i = \hat{z}_j$. Otherwise $V_{i,\cdot} = V_{j,\cdot}$ which implies

$$\|P'_{i,\cdot} - P'_{j,\cdot}\|^2 \leq (\|P'_{i,\cdot} - V_{i,\cdot}\| + \|P'_{j,\cdot} - V_{j,\cdot}\| + \|V_{i,\cdot} - V_{j,\cdot}\|)^2 \leq (p - q)^2 \beta n / k,$$

contradicting with Equation (4.2).

Since \hat{z}_i can only take values in $1, 2, \dots, k$, we conclude $\{\hat{z}_i : i \in C_u\}$ contains only one and different element for all $u \in [k]$. That is, there exists a permutation ρ on $[k]$, such that

$$\hat{z}_i = \rho(u), \forall i \in C_u, \forall u \in [k].$$

Which indicates $\sum_{i \in S^C} \mathbb{I}\{\hat{z}_i \neq \rho(z_i)\} = 0$. Hence

$$\frac{1}{n} \sum_{i \in [n]} \mathbb{I}\{\hat{z}_i \neq \rho(z_i)\} \leq \frac{|S|}{n} \leq \frac{4k^2 c' p}{\beta n (p - q)^2},$$

which holds with probability at least $1 - n^{-2}$. ■

The following lemma on the operator norm of sparse networks is from [12]. In the original statement of Lemma 12 in [12], “with probability $1 - o(1)$ ” is stated. However, its

proof in [12] gives explicit form of the probability that the statement holds, which is at least $1 - n^{-2}$.

Lemma 4.2. [Lemma 12 of [12]] *Suppose M is random symmetric matrix with zero on the diagonal whose entries above the diagonal are independent with the following distribution*

$$M_{i,j} = \begin{cases} 1 - p_{i,j}, & \text{w.p. } p_{i,j}; \\ -p_{i,j}, & \text{w.p. } 1 - p_{i,j}. \end{cases}$$

Let $p \triangleq \max_{i,j} p_{i,j}$ and \tilde{M} be the matrix obtained from M by zeroing out all the rows and columns having more than $20np$ positive entries. Then there exists some constant $c > 0$ such that

$$\|\tilde{M}\|_{\text{op}} \leq c\sqrt{np},$$

holds with probability at least $1 - n^{-2}$.

4.3 Rate-optimal and Computationally Feasible Algorithm

In this section, we propose a polynomial-time algorithm that is also rate-optimal. It consists two parts: an initialization that provides decent network partition and a follow-up refinement that leads to optimal estimation. The refinement step is inspired by the likelihood ratio test for one single node (see Section 4.1). The main idea is as follows. For each node, though we do not know the true labels for the remaining nodes, as long as we have a decent estimation that is close enough to the truth, then the likelihood ratio test proposed in Section 4.1 should also work well. This leads to Algorithm 2 with its theoretical justification presented in Theorem 4.2.

For any $i \in [n]$, define A_{-i} be the matrix after zeroing out the i th row and column of A .

Theorem 4.2. *Under the assumption $z^* \in \mathcal{Z}(n, k, \beta)$ and $\beta^2 n(p-q)^2 / (k^3 p) \rightarrow \infty$, if we use Algorithm 1 as the initial community detection method in Algorithm 2, with high probability,*

Algorithm 2: A Two-stage Algorithm for SBM

Input: Adjacency matrix $A \in \{0, 1\}^{n \times n}$, number of communities k , parameters p, q , initial community detection method \tilde{z}

Output: Partition of the network z

Penalized majority voting:

for $i = 1, 2, \dots, n$ **do**

- 1 Define λ as in Equation (3.6);
- 2 Apply \tilde{z} on A_{-i} to obtain $\tilde{z}_j^{(-i)}$ for all $j \neq i$;
- 3 Define $\hat{z}_j^{(-i)} = \tilde{z}_j^{(-i)}, \forall j \neq i$. Let

$$\hat{z}_i^{(-i)} = \operatorname{argmax}_{u \in [k]} \sum_{j: j \neq i, \hat{z}_j^{(-i)} = u} (A_{i,j} - \lambda).$$

end

Consensus:

- 4 Define $\hat{z}_1 = \hat{z}_1^{(-1)}$. For $i = 2, 3, \dots, n$, define

$$\hat{z}_i = \operatorname{argmax}_{u \in [k]} \left| \{j : \hat{z}_j^{(-1)} = u\} \cap \{j : \hat{z}_j^{(-i)} = \hat{z}_i^{(-i)}\} \right|.$$

we have

$$\ell(\hat{z}, z^*) \leq \begin{cases} \exp(-(1-\eta)nI/2), & k = 2; \\ \exp(-(1-\eta)\beta nI/k), & k \geq 3, \end{cases}$$

for some positive sequence $\eta = o(1)$.

Proof. The proof consists of three parts.

1). Denote \mathcal{F} be the event that $\|T(A) - P\|_{\text{op}} \leq c\sqrt{np}$ for some constant $c_1 > 0$. By Lemma 4.2 we have $\mathbb{P}(\mathcal{F}) \geq 1 - n^{-2}$. Let $\tilde{z}^{(0)}$ be the result after implementing Algorithm 1 on A . From Theorem 4.1, we know $\ell(\tilde{z}^{(0)}, z^*) \leq c_2 k^2 p / (\beta n (p - q)^2)$ holds for some constant

$c_2 > 0$ if event \mathcal{F} holds. Note that

$$\begin{aligned}
\|T(A_{-i}) - P_{-i}\|_{\text{op}} &= \max_{u:\|u\|\leq 1} u^T(T(A_{-i}) - P_{-i})u \\
&= \max_{u:\|u\|\leq 1, u_i=0} u^T(T(A_{-i}) - P_{-i})u \\
&= \max_{u:\|u\|\leq 1, u_i=0} u^T(T(A) - P)u \\
&\leq \|T(A) - P\|_{\text{op}},
\end{aligned}$$

which implies

$$\ell(\tilde{z}^{(-i)}, z^*) \leq \ell(\tilde{z}^{(0)}, z^*) \leq c_2 k^2 p / (\beta n (p - q)^2) + 1/n,$$

holds simultaneously for all $i \in [n]$, by the proof of Theorem 4.1, assuming the event \mathcal{F} holds. The addition $1/n$ term is due to the fact that $\tilde{z}^{(-i)}$ provides no valid value for the i th label.

2). Now we investigate the refinement step for the i th node. Note that we have independence between $\tilde{z}^{(-i)}$ and the data to be used $\{A_{i,j}\}_{j \neq i}$. Without loss of generality, we assume $\|\tilde{z}^{(-i)} - z^*\|_0 = n\ell(\tilde{z}^{(-i)}, z^*)$. Hence, we have

$$\mathbb{P}(\hat{z}_i^{(-i)} \neq z_i^*) \leq \sum_{u \neq z_i^*} \mathbb{P} \left(\sum_{j:j \neq i, \hat{z}_j^{(-i)}=u} (A_{i,j} - \lambda) \geq \sum_{j:j \neq i, \hat{z}_j^{(-i)}=z_i^*} (A_{i,j} - \lambda) \right).$$

If $\ell(\tilde{z}^{(-i)}, z^*) = \eta_1 \beta n / k$ for some $\eta_1 = o(1)$, Lemma 4.3 leads to

$$\begin{aligned}
\mathbb{P}(\hat{z}_i^{(-i)} \neq z_i^*) &\leq \sum_{u \neq z_i^*} \exp \left(-(1 - \eta) \frac{|\{j \neq i : z_j^* = z_i^* \text{ or } u\}| I}{2} \right) \\
&\leq k \exp(-(1 - \eta_2) n_{\min} I),
\end{aligned}$$

where $n_{\min} = \min_{u \neq v} |\{i : z_j^* = u, v\}| / 2$.

3). Now we can combine the above arguments together. If the event \mathcal{F} holds, we have

$$\max_i \ell(\hat{z}^{(-i)}, z^*) \leq \eta_1 \beta n / k$$

holds for some sequence $\eta_1 = o(1)$. The consensus step in Algorithm 2 is essentially to permute all the labels such that

$$\max_i n \left\| \hat{z}^{(-i)} - z^* \right\|_0 \leq \eta_1 \beta n / k.$$

Thus

$$\mathbb{P} \left(\left\{ \hat{z}_i^{(-i)} \neq z_i^* \right\} \cap \mathcal{F} \right) \leq k \exp(-(1 - \eta_2)n_{\min}I), \forall i \in [n],$$

for some $\eta_2 = o(1)$. Define $\eta_3 = \eta_2 + \sqrt{k/(\beta n I)}$. By Markov inequality, if $k \exp(-(1 - \eta_3)n_{\min}I) \geq n^{-3/2}$, then we have

$$\begin{aligned} \mathbb{P}(\|\hat{z} - z^*\|_0 \geq k \exp(-(1 - \eta_3)n_{\min}I)) &\leq \frac{\frac{1}{n} \sum_i \mathbb{P} \left(\left\{ \hat{z}_i^{(-i)} \neq z_i^* \right\} \cap \mathcal{F} \right) + \mathbb{P}(\mathcal{F}^c)}{k \exp(-(1 - \eta_2)n_{\min}I), \forall i \in [n],} \\ &\leq \exp \left(-\sqrt{\frac{\beta n}{k}} \right) + \frac{n^{-2}}{k \exp(-(1 - \eta_3)n_{\min}I)} \\ &\leq \exp \left(-\sqrt{\frac{\beta n}{k}} \right) + \frac{1}{\sqrt{n}} \\ &= o(1). \end{aligned}$$

If $k \exp(-(1 - \eta_3)n_{\min}I) < n^{-3/2}$, we have

$$\begin{aligned} \mathbb{P}(\|\hat{z} - z^*\|_0 \geq k \exp(-(1 - \eta_3)n_{\min}I)) &\leq \mathbb{P}(\|\hat{z} - z^*\|_0 > 0) \\ &\leq \sum_i \mathbb{P} \left(\left\{ \hat{z}_i^{(-i)} \neq z_i^* \right\} \cap \mathcal{F} \right) + \mathbb{P}(\mathcal{F}^c) \\ &\leq 2n^{-3/2} \\ &= o(1). \end{aligned}$$

The proof is complete with

$$n_{\min} \geq \begin{cases} \beta n/k, k \geq 3; \\ n/2, k = 2. \end{cases}$$

■

Lemma 4.3. *Let m_1, m_2 be positive integers. Fix subsets $T_1 \subset [m_1], T_2 \subset [m_2]$ such that $\max\{|T_1|, |T_2|\} = o(\min\{m_1, m_2\})$. Let $U_1, U_2, \dots, U_{m_1}, V_1, V_2, \dots, V_{m_2}$ be mutually independent random variables. Define $X \sim \text{Ber}(q)$ and $Y \sim \text{Ber}(p)$. Let $\{U_i\}_{i \in [m_1] \setminus T_1}$ and $\{V_j\}_{j \in T_2}$ be i.i.d. copies of X ; let $\{U_i\}_{i \in T_1}$ and $\{V_j\}_{j \in [m_2] \setminus T_2}$ be i.i.d. copies of Y . Recall the definition of λ as in Equation (3.6). Under the assumption that $cp \leq q \leq p$ for some positive constant c , we have*

$$\mathbb{P} \left(\sum_{i=1}^{m_1} (U_i - \lambda) \geq \sum_{j=1}^{m_2} (V_j - \lambda) \right) \leq \exp \left(-(1 - \eta) \frac{(m_1 + m_2)I}{2} \right),$$

for some positive sequence $\eta = o(1)$.

Proof. Denote $m'_1 = m_1 - |T_1|$ and $m'_2 = m_2 - |T_2|$. Recall the definition of t as in Equation (3.6). By using Chernoff bound, we have

$$\begin{aligned} & \mathbb{P} \left(\sum_{i=1}^{m_1} (U_i - \lambda) \geq \sum_{j=1}^{m_2} (V_j - \lambda) \right) \\ & \leq e^{-t\lambda(m_1 - m_2)} (\mathbb{E}e^{tX})^{m'_1} (\mathbb{E}e^{tY})^{m_1 - m'_1} (\mathbb{E}e^{-tY})^{m'_2} (\mathbb{E}e^{-tX})^{m_2 - m'_2} \\ & = e^{-t\lambda(m_1 - m_2)} (\mathbb{E}e^{tX})^{m_1} (\mathbb{E}e^{-tY})^{m_2} \left(\frac{\mathbb{E}e^{tY}}{\mathbb{E}e^{tX}} \right)^{m_1 - m'_1} \left(\frac{\mathbb{E}e^{-tX}}{\mathbb{E}e^{-tY}} \right)^{m_2 - m'_2} \\ & = e^{-(m_1 + m_2)I/2} \left(\frac{\mathbb{E}e^{tY}}{\mathbb{E}e^{tX}} \right)^{m_1 - m'_1} \left(\frac{\mathbb{E}e^{-tX}}{\mathbb{E}e^{-tY}} \right)^{m_2 - m'_2}, \end{aligned}$$

where the last equation is due to Proposition 3.4. By the assumption that $q \geq cp$, we have

$$|e^t - 1| = \left| \sqrt{\frac{p(1-q)}{q(1-p)}} - 1 \right| \leq C_1 \frac{p-q}{p},$$

for some constant $C_1 > 0$. Then

$$\frac{\mathbb{E}e^{tY}}{\mathbb{E}e^{tX}} = \frac{pe^t + 1 - p}{qe^{-t} + 1 - q} = 1 + \frac{(p-q)(e^t - 1)}{qe^t + 1 - q} \leq 1 + O\left(\frac{(p-q)^2}{p}\right) \leq \exp\left(O\left(\frac{(p-q)^2}{p}\right)\right).$$

We have the same result for $|e^{-t} - 1|$ and $\mathbb{E}e^{-tX}/\mathbb{E}e^{-tY}$. By Proposition 4.1, due to $\max m_1 - m'_1, m_2 - m'_2 = \max\{|T_1|, |T_2|\} = o(m_1 + m_2)$, we have

$$\mathbb{P}\left(\sum_{i=1}^{m_1}(U_i - \lambda) \geq \sum_{j=1}^{m_2}(V_j - \lambda)\right) \leq \exp(-(1-\eta)(m_1 + m_2)I/2),$$

for some positive sequence $\eta = o(1)$. ■

Proposition 4.1 (Lemma B.1 of [52]). *Let p and q satisfy $\epsilon/n \leq q \leq p \leq 1 - \epsilon$ for any small constant $1 > \epsilon > 0$. We have $I \asymp (p - q)^2/(np)$. In addition if $p = o(1)$, we have $I = (1 + o(1))(\sqrt{p} - \sqrt{q})^2$.*

4.4 Extension

Algorithm 2 is not adaptive, as it requires known p, q . In this section, we present its adaptive counterpart. The main idea is to estimate p, q from the initial partition of the network. On top of that, the proposed method (Algorithm 3) also generalizes Algorithm 2 in another direction, in the sense that it works for inhomogeneous SBM. We are able to provide theoretical guarantees for Algorithm 3 that is very similar to Theorem 4.2. We refer readers to our paper [19] for details.

Algorithm 3: A refinement scheme for community detection

Input: Adjacency matrix $A \in \{0, 1\}^{n \times n}$, number of communities k , initial community detection method \tilde{z}

Output: Community assignment \hat{z} .

Penalized neighbor voting:

1 **for** $i = 1$ **to** n **do**

2 Apply \tilde{z} on A_{-i} to obtain $\tilde{z}_j^{(-i)}$ for all $j \neq i$ and let $\tilde{z}_i^{(-i)} = 0$;

3 Define $\tilde{\mathcal{C}}_u^{(-i)} = \{j : \tilde{z}_j^{(-i)} = u\}$ for all $u \in [k]$; let $\tilde{\mathcal{E}}_u^{(-i)}$ be the set of edges within $\tilde{\mathcal{C}}_u^{(-i)}$, and $\tilde{\mathcal{E}}_{uv}^{(-i)}$ the set of edges between $\tilde{\mathcal{C}}_u^{(-i)}$ and $\tilde{\mathcal{C}}_v^{(-i)}$ when $u \neq v$;

4 Define

$$\hat{B}_{uu}^{(-i)} = \frac{|\tilde{\mathcal{E}}_u^{(-i)}|}{\frac{1}{2}|\tilde{\mathcal{C}}_u^{(-i)}|(|\tilde{\mathcal{C}}_u^{(-i)}| - 1)}, \quad \hat{B}_{u,v}^{(-i)} = \frac{|\tilde{\mathcal{E}}_{uv}^{(-i)}|}{|\tilde{\mathcal{C}}_u^{(-i)}||\tilde{\mathcal{C}}_v^{(-i)}|}, \quad \forall u \neq v \in [k], \quad (4.3)$$

and let

$$\hat{p}_i = \min_{u \in [k]} \hat{B}_{uu}^{(-i)} \quad \text{and} \quad \hat{q}_i = \max_{u \neq v \in [k]} \hat{B}_{uv}^{(-i)}. \quad (4.4)$$

5 Define $\hat{z}^{(-i)} : [n] \rightarrow [k]$ by setting $\hat{z}_j^{(-i)} = \tilde{z}_j^{(-i)}$ for all $j \neq i$ and

$$\hat{z}_i^{(-i)} = \operatorname{argmax}_{u \in [k]} \sum_{\tilde{z}_j^{(-i)} = u} (A_{ij} - \lambda_i) \quad (4.5)$$

where for

$$t_i = \frac{1}{2} \log \frac{\hat{p}_i(1 - \hat{q}_i)}{\hat{q}_i(1 - \hat{p}_i)}, \quad (4.6)$$

we define

$$\lambda_i = -\frac{1}{2t_i} \log \left(\frac{1 - \hat{q}_i}{1 - \hat{p}_i} \right), \quad (4.7)$$

end

Consensus:

6 Define $\hat{z}_1 = \hat{z}_1^{(-1)}$. For $i = 2, \dots, n$, define

$$\hat{z}_i = \operatorname{argmax}_{u \in [k]} \left| \{j : \hat{z}_j^{(-1)} = u\} \cap \{j : \hat{z}_j^{(-i)} = \hat{z}_i^{(-i)}\} \right|. \quad (4.8)$$

Chapter 5

Variational Inference

The Bayesian framework and the variational inference for community detection are considered in [3, 7, 9, 18, 28, 45]. For high dimensional settings, Celisse et al. [9] and Bickel et al. [7] are arguably the first to study the statistical properties of the mean field for SBMs. The authors built an interesting connection between full likelihood and variational likelihood, and then studied the closeness of maximum likelihood and maximum variational likelihood, from which they obtained consistency and asymptotic normality for global parameter estimation. From a personal communication with the authors of Bickel et al. [7], an implication of their results is that the variational method achieves exact community recovery under a strong signal-to-noise (SNR) ratio. Their analysis idea is fascinating, but it is not clear whether it is possible to extend the analysis to other SNR conditions under which exact recovery may never be possible. More importantly, it may not be computationally feasible to maximize the variational likelihood for the SBM, as seen from Theorem 5.1.

In this chapter, we consider the statistical and computational guarantees of the iterative variational inference algorithm for community detection. To the best of our knowledge this provides arguably the first theoretical justification for the iterative algorithm of the mean field variational method in a high-dimensional and complex setting. Though we focus on the problem of community detection in this chapter, we hope the analysis would shed some light on analyzing other models, which may eventually lead to a general framework of understanding the mean field theory.

5.1 Mean Field Variational Inference

We first present the mean field method in a general setting and then consider its application to the community detection problem. Let $\mathbf{p}(x|y)$ be an arbitrary posterior distribution for x , given observation y . Here x can be a vector of latent variables, with coordinates $\{x_i\}$. It may be difficult to compute the posterior $\mathbf{p}(x|y)$ exactly. The variational Bayes ignores the dependence among $\{x_i\}$, by simply taking a product measure $\mathbf{q}(x) = \prod_i \mathbf{q}_i(x_i)$ to approximate it. Usually each $\mathbf{q}_i(x_i)$ is simple and easy to compute. The best approximation is obtained by minimizing the Kullback–Leibler divergence between $\mathbf{q}(x)$ and $\mathbf{p}(x|y)$:

$$\hat{\mathbf{q}}^{\text{MF}} = \underset{\mathbf{q} \in \mathbf{Q}}{\operatorname{argmin}} \operatorname{KL}(\mathbf{q} \parallel \mathbf{p}). \quad (5.1)$$

Despite the fact that every measure \mathbf{q} has a simple product structure, the global minimizer $\hat{\mathbf{q}}^{\text{MF}}$ remains computationally intractable.

To address this issue, an iterative Coordinate Ascent Variational Inference (CAVI) is widely used to approximate the global minimum. It is a greedy algorithm. The value of $\operatorname{KL}(\mathbf{q} \parallel \mathbf{p})$ decreases in each coordinate update:

$$\hat{\mathbf{q}}_i = \min_{\mathbf{q}_i \in \mathbf{Q}_i} \operatorname{KL} \left[\mathbf{q}_i \prod_{j \neq i} \mathbf{q}_j \parallel \mathbf{p} \right], \forall i. \quad (5.2)$$

The coordinate update has an explicit formula

$$\hat{\mathbf{q}}_i(x_i) \propto \exp \left[\mathbb{E}_{\mathbf{q}_{-i}} [\log \mathbf{p}(x_i | x_{-i}, y)] \right], \quad (5.3)$$

where x_{-i} indicates all the coordinates in x except x_i , and the expectation is over $\mathbf{q}_{-i} = \prod_{j \neq i} \mathbf{q}_j(x_j)$. Equation (5.3) is usually easy to compute, which makes CAVI computationally attractive, although CAVI only guarantees to achieve a local minimum.

In summary, the mean field variational inference via CAVI can be represented in the following diagram:

$$\mathbf{p}(x|y) \overset{\text{approx.}}{\longleftarrow} \hat{\mathbf{q}}^{\text{MF}}(x) \overset{\text{approx.}}{\longleftarrow} \hat{\mathbf{q}}^{\text{CAVI}}(x),$$

where $\hat{\mathbf{q}}^{\text{MF}}(x)$, the global minimum, serves mainly as an intermediate step in the mean field methodology. What is implemented in practice to approximate global minimum is an iterative algorithm like CAVI. This motivates us to consider directly the theoretical guarantees of the iterative algorithm in this chapter.

5.2 A Bayesian Framework

The SBM formulated in Chapter 2 can also be written in a matrix form, especially the mean matrix P defined as in Equation (2.1). Let $Z \in \Pi_0$ be the assignment matrix where

$$\Pi_0 = \{\pi \in \{0, 1\}^{n \times k} : \|\pi_{i,\cdot}\|_0 = 1, \forall i \in [n]\}.$$

In each row $\{Z_{i,\cdot}\}_{i=1}^n$ there is only one 1 with all the other coordinates as 0, indicating the assignment of community for the corresponding node. Then P can be equivalently written as $P_{i,j} = Z_{i,\cdot} B Z_{j,\cdot}^T, \forall i < j$, or in a matrix form

$$P_{i,j} = (Z B Z^T)_{i,j}, \forall i < j. \quad (5.4)$$

Consequently, to recover the assignment vector z is equivalent to recover the assignment matrix Z . The equivalence can be seen by observing that there is a bijection r between $z \in [k]^n$ and $Z \in \Pi_0$ which is defined as follows,

$$r(z) = Z, \text{ where } Z_{i,a} = \mathbb{I}\{a = z_i\}, \forall i \in [n], a \in [k]. \quad (5.5)$$

Since they are uniquely determined by each other, in this thesis we may use z directly without explicitly defining $z = r^{-1}(Z)$ (or vice versa) when there is no ambiguity.

Throughout the whole chapter, we assume k , the number of communities, is known. We observe the adjacency matrix A . The global parameters p and q and the community assignment Z are unknown. We can write down the distribution of A as follows:

$$\mathbf{p}(A|Z, p, q) = \prod_{i < j} B_{z_i, z_j}^{A_{i,j}} (1 - B_{z_i, z_j})^{1 - A_{i,j}}, \quad (5.6)$$

with $B = q\mathbf{1}_k\mathbf{1}_k^T + (p - q)I_k$ and $z = r^{-1}(Z)$. We are interested in Bayesian inference for estimating Z , with prior to be given on both p, q and Z .

We assume that $\{z_i\}_{i=1}^n$ have independent categorical (a.k.a. multinomial with size one) priors with hyperparameters $\{\pi_{i,\cdot}^{\text{pri}}\}_{i=1}^n$, where $\sum_{a=1}^k \pi_{i,a}^{\text{pri}} = 1, \forall i \in [n]$. In other words, $\{Z_{i,\cdot}\}_{i=1}^n$ are independently distributed by

$$\mathbb{P}(Z_{i,\cdot} = e_a^T) = \pi_{i,a}^{\text{pri}}, \forall a = 1, 2, \dots, k,$$

where $\{e_a\}_{a=1}^k$ are the coordinate vectors. Here we allow the priors for $Z_{i,\cdot}$ to be different for different i . If additionally $\pi_{i,\cdot} = \pi_{j,\cdot}$ for all $i \neq j$ is assumed, and then this is reduced to the usual case of i.i.d. priors.

Since $\{A_{i,j}\}_{i < j}$ are Bernoulli, it is natural to consider a conjugate Beta prior for p and q . Let $p \sim \text{Beta}(\alpha_p^{\text{pri}}, \beta_p^{\text{pri}})$ and $q \sim \text{Beta}(\alpha_q^{\text{pri}}, \beta_q^{\text{pri}})$. Then the joint distribution is

$$\begin{aligned} \mathbf{p}(A, Z, p, q) &= \left[\prod_i \pi_{i,z_i}^{\text{pri}} \right] \left[\prod_{i < j} B_{z_i, z_j}^{A_{i,j}} (1 - B_{z_i, z_j})^{1 - A_{i,j}} \right] \\ &\times \left[\frac{\Gamma(\alpha_p^{\text{pri}} + \beta_p^{\text{pri}})}{\Gamma(\alpha_p^{\text{pri}})\Gamma(\beta_p^{\text{pri}})} p^{\alpha_p^{\text{pri}} - 1} (1 - p)^{\beta_p^{\text{pri}} - 1} \right] \left[\frac{\Gamma(\alpha_q^{\text{pri}} + \beta_q^{\text{pri}})}{\Gamma(\alpha_q^{\text{pri}})\Gamma(\beta_q^{\text{pri}})} q^{\alpha_q^{\text{pri}} - 1} (1 - q)^{\beta_q^{\text{pri}} - 1} \right]. \end{aligned} \quad (5.7)$$

Our main interest is to infer Z , from the posterior distribution $\mathbf{p}(Z, p, q|A)$. However, the exact calculation of $\mathbf{p}(Z, p, q|A)$ is computationally intractable.

5.3 Mean Field Approximation

Since the posterior distribution $\mathbf{p}(Z, p, q|A)$ is computationally intractable, we apply the mean field approximation to approximate it by a product measure,

$$\mathbf{q}_{\pi, \alpha_p, \beta_p, \alpha_q, \beta_q}(Z, p, q) = \mathbf{q}_{\pi}(Z) \mathbf{q}_{\alpha_p, \beta_p}(p) \mathbf{q}_{\alpha_q, \beta_q}(q)$$

where $\{r^{-1}(Z_{i,\cdot})\}_{i=1}^n$ are independent categorical variables with parameters $\{\pi_{i,\cdot}\}_{i=1}^n$, i.e., $\mathbf{q}_\pi(Z) = \prod_{i=1}^n \mathbf{q}_{\pi_{i,\cdot}}(Z_{i,\cdot})$ with

$$\mathbf{q}_{\pi_{i,\cdot}}(Z_{i,\cdot} = e_a) = \pi_{i,a}, \forall i \in [n], a \in [k],$$

and $\mathbf{q}_{\alpha_p, \beta_p}(p)$ and $\mathbf{q}_{\alpha_q, \beta_q}(q)$ are Beta with parameters $\alpha_p, \beta_p, \alpha_q, \beta_q$ due to conjugacy. See Figure 5.1 for the graphical presentation of $\mathbf{q}_{\pi, \alpha_p, \beta_p, \alpha_q, \beta_q}(Z, p, q)$.

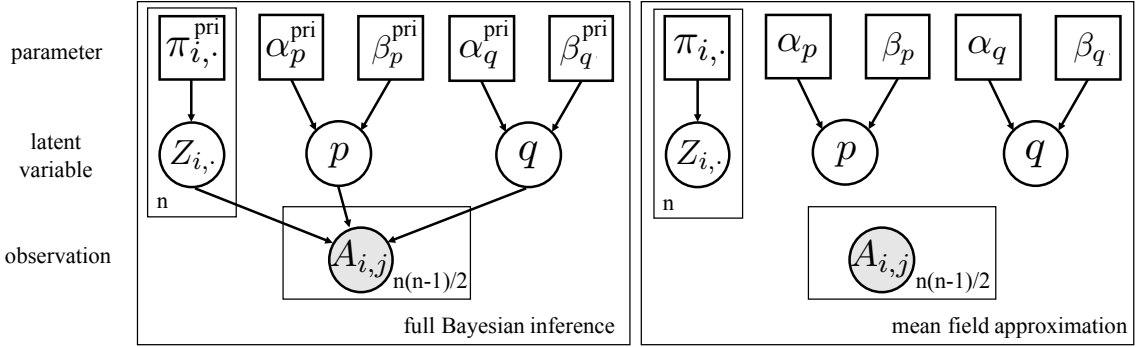


Figure 5.1: Graphical model presentations of full Bayesian inference (*left* panel) and the mean field approximation (*right* panel) for community detection. The edges show the dependence among variables.

Note that the distribution class of \mathbf{q} is fully captured by the parameters $(\pi, \alpha_p, \beta_p, \alpha_q, \beta_q)$, and then the optimization in Equation (5.1) is equivalent to minimize over the parameters as

$$(\hat{\pi}^{\text{MF}}, \hat{\alpha}_p^{\text{MF}}, \hat{\beta}_p^{\text{MF}}, \hat{\alpha}_q^{\text{MF}}, \hat{\beta}_q^{\text{MF}}) = \underset{\substack{\pi \in \Pi_1 \\ \alpha_p, \beta_p, \alpha_q, \beta_q > 0}}{\text{argmin}} \text{KL} \left[\mathbf{q}_{\pi, \alpha_p, \beta_p, \alpha_q, \beta_q}(Z, p, q) \middle| \middle| \mathbf{p}(Z, p, q | A) \right], \quad (5.8)$$

$$\text{where } \Pi_1 = \{\pi \in [0, 1]^{n \times k}, \|\pi_{i,\cdot}\|_1 = 1\}.$$

Here Π_1 can be viewed as a relaxation of Π_0 : it uses an ℓ_1 constraint on each row instead of the ℓ_0 constraint used in Π_0 . The global minimizer $\mathbf{q}_{\hat{\pi}^{\text{MF}}}(Z)$ gives approximate probabilities to classify every node to each community. The optimization in Equation (5.8) can be shown to be equivalent to a more explicit optimization as follows. Recall $\psi(\cdot)$ is the digamma function with $\psi(x) = \frac{d}{dx} [\log \Gamma(x)]$.

Theorem 5.1. *The mean field estimator $(\hat{\pi}^{\text{MF}}, \hat{\alpha}_p^{\text{MF}}, \hat{\beta}_p^{\text{MF}}, \hat{\alpha}_q^{\text{MF}}, \hat{\beta}_q^{\text{MF}})$ defined in Equation*

(5.8) is equivalent to

$$(\hat{\pi}^{MF}, \hat{\alpha}_p^{MF}, \hat{\beta}_p^{MF}, \hat{\alpha}_q^{MF}, \hat{\beta}_q^{MF}) = \underset{\substack{\pi \in \Pi_1 \\ \alpha_p, \beta_p, \alpha_q, \beta_q > 0}}{\operatorname{argmin}} f(\pi, \alpha_p, \beta_p, \alpha_q, \beta_q; A),$$

where

$$\begin{aligned} f(\pi, \alpha_p, \beta_p, \alpha_q, \beta_q; A) &= t \langle A - \lambda \mathbf{1}_n \mathbf{1}_n^T + \lambda I_n, \pi \pi^T \rangle + \frac{1}{2} [\psi(\alpha_q) - \psi(\beta_q)] \|A\|_1 \\ &+ \frac{n}{2} [\psi(\beta_q) - \psi(\alpha_q + \beta_q)] - \sum_{i=1}^n KL \left[\text{Categorical}(\pi_{i,\cdot}) \| \text{Categorical}(\pi_{i,\cdot}^{pri}) \right] \\ &- KL \left[\text{Beta}(\alpha_p, \beta_p) \| \text{Beta}(\alpha_p^{pri}, \beta_p^{pri}) \right] - KL \left[\text{Beta}(\alpha_q, \beta_q) \| \text{Beta}(\alpha_q^{pri}, \beta_q^{pri}) \right], \end{aligned}$$

and

$$t = [[\psi(\alpha_p) - \psi(\beta_p)] - [\psi(\alpha_q) - \psi(\beta_q)]] / 2 \quad (5.9)$$

$$\lambda = [[\psi(\beta_q) - \psi(\alpha_q + \beta_q)] - [\psi(\beta_p) - \psi(\alpha_p + \beta_p)]] / (2t). \quad (5.10)$$

The explicit formulation in Theorem 5.1 is helpful to understand the global minimizer of the mean field method. However, the global minimizer $\hat{\pi}^{MF}$ remains computationally infeasible as the objective function is not convex. Fortunately, there is a practically useful algorithm to approximate it.

5.3.1 Coordinate Ascent Variational Inference

CAVI is possibly the most popular algorithm to approximate the global minimum of the mean field variational Bayes. It is an iterative algorithm. In Equation (5.8), there are latent variables $\{Z_{i,\cdot}\}_{i=1}^n, p, q$. CAVI updates them one by one. Since the distribution class of \mathbf{q} is uniquely determined by the parameters $\{\pi_{i,\cdot}\}_{i=1}^n, \alpha_p, \beta_p, \alpha_q, \beta_q$, equivalently we are updating those parameters iteratively. Theorem 5.2 gives explicit formulas for the coordinate updates.

Theorem 5.2. *Starts with some $\pi, \alpha_p, \beta_p, \alpha_q, \beta_q$, the CAVI update for each coordinate (i.e., Equation (5.2) and Equation (5.3)) has an explicit expression as follows:*

- Update on p :

$$\alpha'_p = \alpha_p^{\text{pri}} + \sum_{i < j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} A_{i,j}, \quad \text{and} \quad \beta'_p = \beta_p^{\text{pri}} + \sum_{i < j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} (1 - A_{i,j}).$$

- Update on q :

$$\alpha'_q = \alpha_q^{\text{pri}} + \sum_{i < j} \sum_{a \neq b} \pi_{i,a} \pi_{j,b} A_{i,j}, \quad \text{and} \quad \beta'_q = \beta_q^{\text{pri}} + \sum_{i < j} \sum_{a \neq b} \pi_{i,a} \pi_{j,b} (1 - A_{i,j}).$$

- Update on $Z_{i,\cdot}, \forall i = 1, 2, \dots, n$:

$$\pi'_{i,a} \propto \pi_{i,a}^{\text{pri}} \exp \left[2t \sum_{j \neq i} \pi_{j,a} (A_{i,j} - \lambda) \right], \quad \forall a = 1, 2, \dots, k,$$

where t and λ are defined in Equation (5.9) and Equation (5.10) respectively, and the normalization satisfies $\sum_{a=1}^k \pi'_{i,a} = 1$.

All coordinate updates in Theorem 5.2 have explicit formulas, which makes CAVI a computationally attractive way to approximate the global optimum $\hat{\mathbf{q}}^{\text{MF}}$ for the community detection problem.

5.3.2 Batch Coordinate Ascent Variational Inference

The Batch Coordinate Ascent Variational Inference (BCAVI) is a batch version of CAVI. The difference lies in that CAVI updates the rows of π sequentially one by one, while BCAVI uses the value of π to update all rows $\{\pi'_{i,\cdot}\}$ according to Theorem 5.2. This makes BCAVI especially suitable for parallel and distributed computing, a nice feature for large scale network analysis.

We define a mapping $h : \Pi_1 \rightarrow \Pi_1$ as follows. For any $\pi \in \Pi_1$, we have

$$[h_{t,\lambda}(\pi)]_{i,a} \propto \pi_{i,a}^{\text{pri}} \exp \left[2t \sum_{j \neq i} \pi_{j,a} (A_{i,j} - \lambda) \right], \quad (5.11)$$

with parameters t and λ . For BCAVI, we update π by $\pi' = h_{t,\lambda}(\pi)$ in each batch iteration, with t, λ defined in Equations (5.14) and (5.15). See Algorithm 4 for the detailed

implementation of BCAVI algorithm.

Algorithm 4: Batch Coordinate Ascent Variational Inference (BCAVI)

Input: Adjacency matrix A , number of communities k , hyperparameters $\pi^{\text{pri}}, \alpha_p^{\text{pri}}, \beta_p^{\text{pri}}, \alpha_q^{\text{pri}}, \beta_q^{\text{pri}}$, initializer $\pi^{(0)}$, number of iterations S .

Output: Mean variational Bayes approximation $\hat{\pi}, \hat{\alpha}_p, \hat{\beta}_p, \hat{\alpha}_q, \hat{\beta}_q$.

for $s = 1, 2, \dots, S$ **do**

1 Update $\alpha_p^{(s)}, \beta_p^{(s)}, \alpha_q^{(s)}, \beta_q^{(s)}$ by

$$\alpha_p^{(s)} = \alpha_p^{\text{pri}} + \sum_{a=1}^k \sum_{i < j} A_{i,j} \pi_{i,a}^{(s-1)} \pi_{j,a}^{(s-1)}, \beta_p^{(s)} = \beta_p^{\text{pri}} + \sum_{a=1}^k \sum_{i < j} (1 - A_{i,j}) \pi_{i,a}^{(s-1)} \pi_{j,a}^{(s-1)}, \quad (5.12)$$

$$\alpha_q^{(s)} = \alpha_q^{\text{pri}} + \sum_{a \neq b} \sum_{i < j} A_{i,j} \pi_{i,a}^{(s-1)} \pi_{j,b}^{(s-1)}, \beta_q^{(s)} = \beta_q^{\text{pri}} + \sum_{a \neq b} \sum_{i < j} (1 - A_{i,j}) \pi_{i,a}^{(s-1)} \pi_{j,b}^{(s-1)}. \quad (5.13)$$

2 Define

$$t^{(s)} = \frac{1}{2} \left[\left[\psi(\alpha_p^{(s)}) - \psi(\beta_p^{(s)}) \right] - \left[\psi(\alpha_q^{(s)}) - \psi(\beta_q^{(s)}) \right] \right] \quad (5.14)$$

$$\lambda^{(s)} = \frac{1}{2t^{(s)}} \left[\left[\psi(\beta_q^{(s)}) - \psi(\alpha_q^{(s)} + \beta_q^{(s)}) \right] - \left[\psi(\beta_p^{(s)}) - \psi(\alpha_p^{(s)} + \beta_p^{(s)}) \right] \right], \quad (5.15)$$

where $\psi(\cdot)$ is the digamma function. Then update $\pi^{(s)}$ with

$$\pi^{(s)} = h_{t^{(s)}, \lambda^{(s)}}(\pi^{(s-1)}),$$

where the mapping $h(\cdot)$ is defined as in Equation (5.11).

end

3 We have $\hat{\pi} = \pi^{(S)}, \hat{\alpha}_p = \alpha_p^{(S)}, \hat{\beta}_p = \beta_p^{(S)}, \hat{\alpha}_q = \alpha_q^{(S)}, \hat{\beta}_q = \beta_q^{(S)}$.

Remark 1. The definitions of $t^{(s)}$ and $\lambda^{(s)}$ in Equations (5.14) and (5.15) involve the digamma function, which costs a non-negligible computational resources each time called. Note that we have $\psi(x) \in (\log(x - \frac{1}{2}), \log x)$ for all $x > 1/2$. For the computational purpose, we propose to use the logarithmic function instead of digamma function in Algorithm 4, i.e., Equations (5.14) and (5.15) are replaced by

$$t^{(s)} = \frac{1}{2} \log \frac{\alpha_p^{(s)} \beta_q^{(s)}}{\beta_p^{(s)} \alpha_q^{(s)}}, \quad \text{and } \lambda^{(s)} = \frac{1}{2t^{(s)}} \log \frac{\beta_q^{(s)} (\alpha_p^{(s)} + \beta_p^{(s)})}{(\alpha_q^{(s)} + \beta_q^{(s)}) \beta_p^{(s)}}. \quad (5.16)$$

Later we show that $\alpha_p^{(s)}, \beta_p^{(s)}, \alpha_q^{(s)}, \beta_q^{(s)}$ are all at least in the order of np , which goes to infinity, and thus the error caused by using the logarithmic function to replace the digamma

function is negligible. All theoretical guarantees obtained in Section 5.4 for Algorithm 4 (i.e., Theorem 5.3, Theorem 5.4) still hold if we use Equation (5.16) to replace Equations (5.14) and (5.15).

5.4 Theoretical Justifications

In this section, we establish theoretical justifications for BCAVI for community detection under the Stochastic Block Model. Though Z , p and q are all unknown, the main interest of community detection is on the recovery of the assignment matrix Z , while p and q are nuisance parameters. As a result, our main focus is on developing convergence rate of BCAVI for π .

5.4.1 Loss Function

We use ℓ_1 norm to measure the performance of recovering Z . Then for any $Z, Z^* \in \Pi_1$, the loss function is defined as

$$\ell(Z, Z^*) = \frac{1}{n} \min_{\rho} \|Z - \rho \circ Z^*\|_1 = \frac{1}{n} \min_{\rho} \sum_{i,a} |Z_{i,a} - Z_{i,\rho(a)}^*|, \quad (5.17)$$

where the minimization is over all permutations on $[k]$ to avoid an identifiability issue of labels.

There are a few reasons for the choice of the ℓ_1 norm. When both $Z, Z' \in \Pi_0$, the ℓ_1 distance between Z and Z' is equal to the ℓ_0 norm, i.e., the Hamming distance between the corresponding assignment vectors $r^{-1}(Z)$ and $r^{-1}(Z')$, which matches with the distance used in the previous chapters. Despite a little abuse of notation, we use the same $\ell(\cdot, \cdot)$ notation here.

5.4.2 Ground Truth

We use the superscript asterisk (*) to indicate the ground truth. The ground truth of connectivity matrix B^* is

$$B^* = q^* \mathbf{1}_k \mathbf{1}_k^T + (p^* - q^*) I_k,$$

where p^* is the within community connection probability and q^* is the between community connection probability. Throughout the chapter, we assume $p^* > q^*$ such that the network satisfies the so-called “assortative” property, with the within-community connectivity probability larger than the between-community connectivity probability.

We further assume the network is generated by the true assignment matrix Z^* in the sense that $P_{i,j} = (Z^* B^* Z^{*T})_{i,j}$ for all $i \neq j$. We are interested in deriving a statistical guarantee of $\ell(\hat{\pi}^{(s)}, Z^*)$. Throughout this section we consider cases $Z^* \in \Pi_0$ or $Z^* \in \Pi_0^{(\rho, \rho')}$, where $\Pi_0^{(\rho, \rho')}$ is defined to be a subset of Π_0 with all the community sizes bounded between $\rho n/k$ and $\rho' n/k$. That is,

$$\Pi_0^{(\rho, \rho')} = \{\pi \in \Pi_0 : \rho n/k \leq |\{i \in [n] : \pi_{i,a} = 1\}| \leq \rho' n/k, \forall a \in [k]\}.$$

It is worth mentioning that ρ, ρ' are not necessarily constants. We allow the community sizes not to be of the same order in the theoretical analysis.

5.4.3 Guarantees

In Theorem 5.3, we present theoretic guarantees of the convergence rate of BCAVI when initialized properly. Define

$$w = \max_{i \in [n]} \max_{a, b \in [k]} \pi_{i,a}^{\text{pri}} / \pi_{i,b}^{\text{pri}}, \text{ and } \bar{n}_{\min} = \min_{a \neq b} [n_a + n_b] / 2.$$

When $w = 1$, the priors for $\{r^{-1}(Z_{i,\cdot})\}_{i=1}^n$ are Categorical with parameter $(1/k, 1/k, \dots, 1/k)$ and $\bar{n}_{\min} = n/2$ when there exist only two communities. The following quantity I plays a

key role in the minimax theory [52]

$$I = -2 \log \left[\sqrt{p^* q^*} + \sqrt{(1-p^*)(1-q^*)} \right],$$

which is the Rényi divergence of order 1/2 between two Bernoulli distributions: $\text{Ber}(p^*)$ and $\text{Ber}(q^*)$. The proof of Theorem 5.3 is deferred to Section 5.5.3.

Theorem 5.3. *Let $Z^* \in \Pi_0$. Let $0 < c_0 < 1$ be any constant. Assume $0 < c_0 p^* < q^* < p^* = o_n(1)$,*

$$nI/[wk[n/\bar{n}_{\min}]^2] \rightarrow \infty, \text{ and } \alpha_p^{\text{pri}}, \beta_p^{\text{pri}}, \alpha_q^{\text{pri}}, \beta_q^{\text{pri}} = o_n((p^* - q^*)n^2/k). \quad (5.18)$$

Under the assumption that the initializer $\pi^{(0)}$ satisfies $\ell(\pi^{(0)}, Z^) \leq c_{\text{init}} \bar{n}_{\min}$ for some sufficiently small constant c_{init} with probability at least $1 - \epsilon$, there exist some constant $c > 0$ and some $\eta = o_n(1)$ such that in each iteration for the BCAVI algorithm, we have*

$$\ell(\pi^{(s+1)}, Z^*) \leq \exp(-(1-\eta)\bar{n}_{\min}I) + \frac{\ell(\pi^{(s)}, Z^*)}{\sqrt{nI/[wk[n/\bar{n}_{\min}]^2]}}, \forall s \geq 0,$$

holds uniformly with probability at least $1 - \exp[-(\bar{n}_{\min}I)^{\frac{1}{2}}] - n^{-c} - \epsilon$.

Theorem 5.3 establishes a linear convergence rate for BCAVI algorithm. The coefficient $[nI/[wk[n/\bar{n}_{\min}]^2]]^{-1/2}$ is independence of s , and goes to 0 when n grows. The following theorem is an immediate consequence of Theorem 5.3.

Theorem 5.4. *Under the same condition as in Theorem 5.3, for any*

$$s \geq s_0 \triangleq [nI/k]/\log[nI/[wk[n/\bar{n}_{\min}]^2]],$$

we have

$$\ell(\hat{\pi}^{(s)}, Z^*) \leq \exp(-(1-2\eta)\bar{n}_{\min}I) \leq \begin{cases} \exp(-(1-o(1))\rho nI/k), & k \geq 3; \\ \exp(-(1-o(1))nI/2), & k = 2, \end{cases}$$

with probability at least $1 - \exp[-(\bar{n}_{\min}I)^{\frac{1}{2}}] - n^{-c} - \epsilon$.

Theorem 5.4 shows that BCAVI provably attains the statistical optimality from the minimax lower bound in Theorem 3.1 after at most s_0 iterations. When the network is sparse, i.e., p^* and q^* are at most in an order of $(\log n)/n$, the quantity s_0 can be shown to be $o(\log n)$, and then BCAVI converges to be minimax rate within $\log n$ iterations. When the network is dense, i.e., p^* and q^* are far bigger than $(\log n)/n$, $\log n$ iterations are not enough to attain the minimax rate. However, $n\ell(\pi^{(s)}, Z^*) = o(n^{-a})$ for any $a > 0$ when $s \geq \log n$, and thus all the nodes can be correctly clustered with high probability by clustering each node to a community with the highest assignment probability. Therefore, it is enough to pick the number of iterations to be $\log n$ in implementing BCAVI.

To help understand Theorem 5.3, we add a remark on conditions on model parameters and priors, and a remark on initialization.

Remark 1 (Conditions on model parameters and priors). The community sizes are not necessarily of the same order in Theorem 5.3. If we further assume ρ, ρ' are constants, and the prior $\pi_{i,a}^{\text{pri}} \asymp 1/k, \forall i \in [n], a \in [k]$ (for example, uniform prior), and then the first condition in Equation (5.18) is equivalent to

$$nI/k^3 \rightarrow \infty,$$

noting that $n/\bar{n}_{\min} \asymp k$ and $w \asymp 1$. It is comparable to the condition in Theorem 4.2.

Under the assumption $nI/k^3 \rightarrow \infty$, since we have $I \asymp (p^* - q^*)^2/p^*$, it can be shown that p^*, q^* are far bigger than n^{-1} , and then the second part of Equation (5.18) can also be easily satisfied. For instance, we can simply set $\alpha_p^{\text{pri}}, \beta_p^{\text{pri}}, \alpha_q^{\text{pri}}, \beta_q^{\text{pri}}$ all equals to 1, i.e., consider non-informative priors.

Remark 2 (Initialization). The requirement on the initializers for BCAVI in Theorem 5.3 is relatively weak. When k is a constant and the community sizes are of the same order, the condition needed is $\ell(\pi^{(0)}, Z^*) \leq c$ for some small constant c . Many existing methodologies in community detection literature can be used. One popular choice is spectral clustering, or our low-rank based spectral clustering proposed in Algorithm 1. They have a mis-clustering

error bound as $\mathcal{O}(k^2/I)$. From Equation (5.18), the error is $o(\bar{n}_{\min})$, and then the condition that Theorem 5.3 requires for initialization is satisfied. The semidefinite programming (SDP), another popular method for community detection, also enjoys satisfactory theoretical guarantees [16, 23], and is suitable as an initializer.

5.5 Proofs

5.5.1 Proof of Theorem 5.1

From Equation (5.8), by some algebra we have

$$(\hat{\pi}^{\text{MF}}, \hat{\alpha}_p^{\text{MF}}, \hat{\beta}_p^{\text{MF}}, \hat{\alpha}_q^{\text{MF}}, \hat{\beta}_q^{\text{MF}}) = \underset{\substack{\pi \in \Pi_1 \\ \alpha_p, \beta_p, \alpha_q, \beta_q > 0}}{\text{argmin}} \mathbb{E}_{\mathbf{q}}[\log \mathbf{p}(A|Z, p, q)] - \text{KL}(\mathbf{q}(Z, p, q) \parallel \mathbf{p}(Z, p, q)), \quad (5.19)$$

where we use \mathbf{q} instead of $\mathbf{q}_{\pi, \alpha_p, \beta_p, \alpha_q, \beta_q}$ for simplicity. From the conditional distribution in Equation (5.6), the log-likelihood function can be simplified as

$$\log \mathbf{p}(A|Z, p, q) = \sum_{a,b} \sum_{i < j} Z_{ia} Z_{jb} \left[A_{i,j} \log \frac{B_{ab}}{1 - B_{ab}} + \log(1 - B_{ab}) \right].$$

Due to the independence of Z and p, q under \mathbf{q} , we have

$$\begin{aligned} \mathbb{E}_{\mathbf{q}}[\log \mathbf{p}(A|Z, p, q)] &= \mathbb{E}_{\mathbf{q}(p,q)} \left[\mathbb{E}_{\mathbf{q}(Z)} \left[\sum_{a,b} \sum_{i < j} Z_{i,a} Z_{j,b} \left[A_{i,j} \log \frac{B_{ab}}{1 - B_{ab}} + \log(1 - B_{ab}) \right] \right] \right] \\ &= \mathbb{E}_{\mathbf{q}(p,q)} \left[\sum_{a,b} \sum_{i < j} \pi_{i,a} \pi_{j,b} \left[A_{i,j} \log \frac{B_{ab}}{1 - B_{ab}} + \log(1 - B_{ab}) \right] \right]. \end{aligned}$$

Since $B_{a,a} = p, \forall a \in [k]$ and $B_{a,b} = q, \forall a \neq b$, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{q}}[\log \mathbf{p}(A|Z, p, q)] &= \mathbb{E}_{\mathbf{q}(p,q)} \left[\sum_a \sum_{i < j} \pi_{i,a} \pi_{j,a} \left[A_{i,j} \log \frac{p(1-q)}{q(1-p)} + \log \frac{1-p}{1-q} \right] \right] \\ &\quad + \mathbb{E}_{\mathbf{q}(p,q)} \left[\sum_{a,b} \sum_{i < j} \pi_{i,a} \pi_{j,b} \left[A_{i,j} \log \frac{q}{1-q} + \log(1-q) \right] \right]. \end{aligned} \quad (5.20)$$

By properties of Beta distribution, we obtain

$$\begin{aligned}\mathbb{E}_{\mathbf{q}(p,q)} \log \frac{p(1-q)}{q(1-p)} &= \mathbb{E}_{\mathbf{q}(p)} [\log p - \log(1-p)] - \mathbb{E}_{\mathbf{q}(q)} [\log q - \log(1-q)] \\ &= [\psi(\alpha_p) - \psi(\beta_p)] - [\psi(\alpha_q) - \psi(\beta_q)],\end{aligned}$$

and

$$\begin{aligned}\mathbb{E}_{\mathbf{q}(p,q)} \log \frac{1-q}{1-p} &= \mathbb{E}_{\mathbf{q}(q)} \log(1-q) - \mathbb{E}_{\mathbf{q}(p)} \log(1-p) \\ &= [\psi(\beta_q) - \psi(\alpha_q + \beta_q)] - [\psi(\beta_p) - \psi(\alpha_p + \beta_p)].\end{aligned}$$

This leads to

$$\begin{aligned}\mathbb{E}_{\mathbf{q}(p,q)} \left[\sum_a \sum_{i < j} \pi_{i,a} \pi_{j,a} \left[A_{i,j} \log \frac{p(1-q)}{q(1-p)} + \log \frac{1-p}{1-q} \right] \right] &= 2t \left[\sum_a \sum_{i < j} \pi_{i,a} \pi_{j,a} (A_{i,j} - \lambda) \right] \\ &= t \langle A - \lambda \mathbf{1}_n \mathbf{1}_n^T + \lambda I_n, \pi \pi^T \rangle.\end{aligned}\tag{5.21}$$

Similarly we can obtain

$$\begin{aligned}\mathbb{E}_{\mathbf{q}(p,q)} \left[\sum_{a,b} \sum_{i < j} \pi_{i,a} \pi_{j,b} \left[A_{i,j} \log \frac{q}{1-q} + \log(1-q) \right] \right] &= \left[\mathbb{E}_{\mathbf{q}(q)} \log \frac{q}{1-q} \right] \sum_{i < j} A_{i,j} \sum_{a,b} \pi_{i,a} \pi_{j,b} + \left[\mathbb{E}_{\mathbf{q}(q)} \log(1-q) \right] \sum_{i < j} \sum_{a,b} \pi_{i,a} \pi_{j,b} \\ &= \frac{1}{2} [\psi(\alpha_q) - \psi(\beta_q)] \|A\|_1 + \frac{n}{2} [\psi(\beta_q) - \psi(\alpha_q + \beta_q)],\end{aligned}\tag{5.22}$$

where we use the fact that $\|\pi_{i,\cdot}\|_1 = 1, \forall i \in [n]$. Now consider the Kullback–Leibler divergence between $\mathbf{q}(Z, p, q)$ and $\mathbf{p}(Z, p, q)$. Due to the independence of p, q and $\{Z_{i,\cdot}\}_{i=1}^n$

in both distributions, we have

$$\begin{aligned}
\text{KL}(\mathbf{q}(Z, p, q) \| \mathbf{p}(Z, p, q)) &= \text{KL}(\mathbf{q}(Z) \| \mathbf{p}(Z)) + \text{KL}(\mathbf{q}(p) \| \mathbf{p}(p)) + \text{KL}(\mathbf{q}(q) \| \mathbf{p}(q)) \quad (5.23) \\
&= \sum_{i=1}^n \text{KL} \left[\text{Categorical}(\pi_{i,\cdot}) \| \text{Categorical}(\pi_i^{\text{pri}}) \right] \\
&\quad + \text{KL} \left[\text{Beta}(\alpha_p, \beta_p) \| \text{Beta}(\alpha_p^{\text{pri}}, \beta_p^{\text{pri}}) \right] + \text{KL} \left[\text{Beta}(\alpha_q, \beta_q) \| \text{Beta}(\alpha_q^{\text{pri}}, \beta_q^{\text{pri}}) \right].
\end{aligned}$$

By Equations (5.19) - (5.23), we conclude with the desired result.

5.5.2 Proof of Theorem 5.2

Note that

$$B_{z_i, z_j} = \left[\sum_{a=1}^k Z_{i,a} Z_{j,a} \right] p + \left[\sum_{a \neq b} Z_{i,a} Z_{j,b} \right] q.$$

We rewrite the joint distribution $\mathbf{p}(p, q, z, A)$ in Equation (5.7) as follows,

$$\begin{aligned}
\mathbf{p}(p, q, Z, A) & \quad (5.24) \\
&= \left[\prod_{i=1}^n \pi_{i, z_i}^{\text{pri}} \right] \left[\prod_{i < j} [p^{A_{i,j}} (1-p)^{1-A_{i,j}}]^{\sum_{a=1}^k Z_{i,a} Z_{j,a}} \right] \left[\prod_{i < j} [q^{A_{i,j}} (1-q)^{1-A_{i,j}}]^{\sum_{a \neq b} Z_{i,a} Z_{j,b}} \right] \\
&\quad \times \left[\frac{\Gamma(\alpha_p^{\text{pri}} + \beta_p^{\text{pri}})}{\Gamma(\alpha_p^{\text{pri}}) \Gamma(\beta_p^{\text{pri}})} p^{\alpha_p^{\text{pri}} - 1} (1-p)^{\beta_p^{\text{pri}} - 1} \right] \left[\frac{\Gamma(\alpha_q^{\text{pri}} + \beta_q^{\text{pri}})}{\Gamma(\alpha_q^{\text{pri}}) \Gamma(\beta_q^{\text{pri}})} q^{\alpha_q^{\text{pri}} - 1} (1-q)^{\beta_q^{\text{pri}} - 1} \right].
\end{aligned}$$

Updates on p and q From Equation (5.24), p has conditional probability as

$$\mathbf{p}(p|q, Z, A) \propto \left[\prod_{i < j} [p^{A_{i,j}} (1-p)^{1-A_{i,j}}]^{\sum_{a=1}^k Z_{i,a} Z_{j,a}} \right] \left[\frac{\Gamma(\alpha_p^{\text{pri}} + \beta_p^{\text{pri}})}{\Gamma(\alpha_p^{\text{pri}}) \Gamma(\beta_p^{\text{pri}})} p^{\alpha_p^{\text{pri}} - 1} (1-p)^{\beta_p^{\text{pri}} - 1} \right].$$

Then the CAVI update in Equation (5.3) leads to

$$\begin{aligned}
\hat{\mathbf{q}}(p) &\propto \exp \left[\mathbb{E}_{\mathbf{q}(q,Z)} \log \mathbf{p}(p|q, Z, A) \right] \\
&\propto \exp \left[\mathbb{E}_{\mathbf{q}(Z)} \sum_{i < j} \sum_{a=1}^k Z_{i,a} Z_{j,a} \log [p^{A_{i,j}} (1-p)^{1-A_{i,j}}] \right] \left[\frac{\Gamma(\alpha_p^{\text{pri}} + \beta_p^{\text{pri}})}{\Gamma(\alpha_p^{\text{pri}}) \Gamma(\beta_p^{\text{pri}})} p^{\alpha_p^{\text{pri}} - 1} (1-p)^{\beta_p^{\text{pri}} - 1} \right] \\
&= \exp \left[\sum_{i < j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} \log [p^{A_{i,j}} (1-p)^{1-A_{i,j}}] \right] \left[\frac{\Gamma(\alpha_p^{\text{pri}} + \beta_p^{\text{pri}})}{\Gamma(\alpha_p^{\text{pri}}) \Gamma(\beta_p^{\text{pri}})} p^{\alpha_p^{\text{pri}} - 1} (1-p)^{\beta_p^{\text{pri}} - 1} \right].
\end{aligned}$$

It can be written as

$$\hat{\mathbf{q}}(p) \propto \left[p^{\sum_{i < j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} A_{i,j}} (1-p)^{\sum_{i < j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} (1-A_{i,j})} \right] \left[\frac{\Gamma(\alpha_p^{\text{pri}} + \beta_p^{\text{pri}})}{\Gamma(\alpha_p^{\text{pri}}) \Gamma(\beta_p^{\text{pri}})} p^{\alpha_p^{\text{pri}} - 1} (1-p)^{\beta_p^{\text{pri}} - 1} \right].$$

The distribution of p is still Beta $p \sim \text{Beta}(\alpha'_p, \beta'_p)$, with

$$\alpha'_p = \alpha_p^{\text{pri}} + \sum_{i < j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} A_{i,j}, \text{ and } \beta'_p = \beta_p^{\text{pri}} + \sum_{i < j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} (1 - A_{i,j}).$$

Similar analysis on q yields updates on α'_q and β'_q . Hence, its proof is omitted.

Updates on $\{Z_{i,\cdot}\}_{i=1}^n$ From Equation (5.24), the conditional distribution on $Z_{i,\cdot}$ is

$$\mathbf{p}(Z_{i,\cdot} | Z_{-i,\cdot}, p, q, A) \propto \pi_{i,z_i}^{\text{pri}} \left[\prod_{j \neq i} B_{z_i, z_j}^{A_{i,j}} (1 - B_{z_i, z_j})^{1-A_{i,j}} \right].$$

Consequently, up to a constant not depending on i , we have

$$\begin{aligned}
&\log \mathbb{P}(Z_{i,a} = 1 | Z_{-i,\cdot}, p, q, A) \\
&= \log \pi_{i,a}^{\text{pri}} + \log \left[\sum_{j \neq i} Z_{j,a} \left[A_{i,j} \log \frac{p}{1-p} + \log(1-p) \right] + \sum_{j \neq i} \sum_{b \neq a} Z_{j,b} \left[A_{i,j} \log \frac{q}{1-q} + \log(1-q) \right] \right] \\
&= \log \pi_{i,a}^{\text{pri}} + \log \left[\sum_{j \neq i} Z_{j,a} \left[A_{i,j} \log \frac{p(1-q)}{q(1-p)} - \log \frac{1-q}{1-p} \right] + \sum_{j \neq i} \left[A_{i,j} \log \frac{q}{1-q} + \log(1-q) \right] \right].
\end{aligned}$$

Then the CAVI update from Equation (5.3) leads to

$$\begin{aligned}
\pi'_{i,a} &= \hat{\mathbf{q}}_{Z_{i,\cdot}}(Z_{i,a} = 1) \\
&\propto \exp \left[\mathbb{E}_{\mathbf{q}(p,q,z_{-i})} \log \mathbb{P}(Z_{i,a} = 1 | Z_{-i,\cdot}, p, q, A) \right] \\
&= \exp \left[\mathbb{E}_{\mathbf{q}(p)} \mathbb{E}_{\mathbf{q}(q)} \mathbb{E}_{\mathbf{q}(Z_{-i,\cdot})} \log \mathbb{P}(Z_{i,a} = 1 | Z_{-i,\cdot}, p, q, A) \right] \\
&\propto \pi_{i,a}^{\text{pri}} \exp \left[\mathbb{E}_{\mathbf{q}(p)} \mathbb{E}_{\mathbf{q}(q)} \sum_{j \neq i} \pi_{j,a} \left[A_{i,j} \log \frac{p(1-q)}{q(1-p)} - \log \frac{1-q}{1-p} \right] \right], \quad (5.25)
\end{aligned}$$

where we use the property that p, q, Z are all independent of each other under \mathbf{q} . Recall that $p \sim \text{Beta}(\alpha_p, \beta_p)$ and $q \sim \text{Beta}(\alpha_q, \beta_q)$. It can be shown that

$$\mathbb{E}_{\mathbf{q}(p)} \log \frac{p}{1-p} = \psi(\alpha_p) - \psi(\beta_p), \text{ and } \mathbb{E}_{\mathbf{q}(p)} \log(1-p) = \psi(\beta_p) - \psi(\alpha_p + \beta_p),$$

where $\psi(\cdot)$ is digamma function. Similar results hold for $\mathbb{E}_{\mathbf{q}(q)} \log(q/(1-q))$ and $\mathbb{E}_{\mathbf{q}(q)} \log(1-q)$. Plug in these expectations to Equation (5.25), we have

$$\pi'_{i,a} \propto \pi_{i,a}^{\text{pri}} \exp \left[2t \sum_{j \neq i} \pi_{j,a} (A_{i,j} - \lambda) \right].$$

5.5.3 Proof of Theorem 5.3

Theorem 5.3 gives a theoretical justification for all iterations in the BCAVI algorithm. Due to the limit of pages, in this section we assume $n\ell(\pi^{(0)}, Z^*) = o(\bar{n}_{\min})$. The proof of the case $n\ell(\pi^{(0)}, Z^*)$ in a constant order of \bar{n}_{\min} is essentially the same with slight modification and thus omitted here.

To prove the theorem, it is sufficient if we are able to show the loss $\ell(\cdot, Z^*)$ decreases in a desired way for one BCAVI iteration, when the community assignment is in an appropriate neighborhood of the truth. Let $\gamma = o(1)$ be any sequence that goes to zero when n grows. Define t^* and λ^* as the true counterparts of t and λ , by

$$t^* = \frac{1}{2} \log \frac{p^*(1-q^*)}{q^*(1-p^*)}, \text{ and } \lambda^* = \frac{1}{2t^*} \log \frac{1-q^*}{1-p^*}.$$

The proof of Theorem 5.3 involves three parts as follows.

Part One: One Iteration. Consider any $\pi \in \Pi_1$ such that $\|\pi - Z^*\|_1 \leq \gamma \bar{n}_{\min}$. Let η' be any sequence such that $\eta' = o(1)$. Consider any t and λ with $|t - t^*| \leq \eta'(p^* - q^*)/p^*$ and $|\lambda - \lambda^*| \leq \eta'(p^* - q^*)$. We define \mathcal{F} to be the event, that after applying the mapping $h_{t,\lambda}(\cdot)$, there exists some $\eta = o(1)$ such that

$$\|h_{t,\lambda}(\pi) - Z^*\|_1 \leq n \exp(-(1 - \eta)\bar{n}_{\min}I) + \frac{\|\pi - Z^*\|_1}{\sqrt{nI/[wk[n/\bar{n}_{\min}]^2]}},$$

holds uniformly over all the eligible π, t and λ . We have

$$\mathbb{P}(\mathcal{F}) \geq 1 - \exp[-(\bar{n}_{\min}I)^{\frac{1}{2}}] - n^{-r},$$

for some constant $r > 0$. We defer its proof to the later part of this section.

Part Two: Consistency of Model Parameters. Consider any $\pi \in \Pi_1$ such that $\|\pi - Z^*\|_1 \leq \gamma \bar{n}_{\min}$. Define

$$\alpha_p = \alpha_p^{\text{pri}} + \sum_{a=1}^k \sum_{i < j} A_{i,j} \pi_{i,a} \pi_{j,a}, \quad \beta_p = \beta_p^{\text{pri}} + \sum_{a=1}^k \sum_{i < j} (1 - A_{i,j}) \pi_{i,a} \pi_{j,a}, \quad (5.26)$$

and

$$\alpha_q = \alpha_q^{\text{pri}} + \sum_{a \neq b} \sum_{i < j} A_{i,j} \pi_{i,a} \pi_{j,b}, \quad \beta_q = \beta_q^{\text{pri}} + \sum_{a \neq b} \sum_{i < j} (1 - A_{i,j}) \pi_{i,a} \pi_{j,b}, \quad (5.27)$$

and consequently,

$$t = \frac{1}{2} [[\psi(\alpha_p) - \psi(\beta_p)] - [\psi(\alpha_q) - \psi(\beta_q)]] \quad (5.28)$$

$$\lambda = \frac{1}{2t} [[\psi(\beta_q) - \psi(\alpha_q + \beta_q)] - [\psi(\beta_p) - \psi(\alpha_p + \beta_p)]]. \quad (5.29)$$

From Lemma 5.1, we have a concentration of t, λ towards t^*, λ^* . That is, there exists some

$\eta' = o(1)$, such that with probability at least $1 - e^{35^{-n}}$, the following inequalities hold

$$|t - t^*| \leq \eta'(p^* - q^*)/p^*, \text{ and } |\lambda - \lambda^*| \leq \eta'(p^* - q^*),$$

uniformly over all the eligible π .

Part Three: Multiple Iterations. Consider any $\pi \in \Pi_1$ such that $\|\pi - Z^*\|_1 \leq \gamma \bar{n}_{\min}$. Define $\alpha_p, \beta_p, \alpha_q, \beta_q, t, \lambda$ as Equations (5.26) - (5.29). A combination of results from *Part One* and *Part Two* immediately implies that

$$\|h_{t,\lambda}(\pi) - Z^*\|_1 \leq n \exp(-(1 - \eta)\bar{n}_{\min}I) + \frac{\|\pi - Z^*\|_1}{\sqrt{nI/[wk[n/\bar{n}_{\min}]^2]}}, \quad (5.30)$$

holds uniformly over all the eligible π with probability at least $1 - \exp[-(\bar{n}_{\min}I)^{\frac{1}{2}}] - n^{-r}$.

This is sufficient to show Theorem 5.3.

The only thing left to be proved, the most critical part towards the proof of Theorem 5.3, is the claim we made in *Part One*. We are going to prove the claim as follow.

Proof Sketch of Part One. The error associated with the $[h_{t,\lambda}(\pi)]_{i,\cdot}$ is a function of π and $A_{i,\cdot}$. It can be decomposed into a summation of two terms, one only involves the ground truth Z^* and the other involves the deviation $\pi - Z^*$. That is,

$$\|[h_{t,\lambda}(\pi)]_{i,\cdot} - Z_{i,\cdot}^*\|_1 \leq f_{i,1}(Z^*, A_{i,\cdot}) + f_{i,2}(\pi - Z^*, A_{i,\cdot}).$$

Consequently,

$$\|h_{t,\lambda}(\pi) - Z^*\|_1 \leq \underbrace{\sum_{i=1}^n f_{i,1}(Z^*, A_{i,\cdot})}_{\text{involves } Z^*} + \underbrace{\sum_{i=1}^n f_{i,2}(\pi - Z^*, A_{i,\cdot})}_{\text{involves } \pi - Z^*}. \quad (5.31)$$

With a proper choice of $f_{\cdot,1}$ and $f_{\cdot,2}$, the first term on the RHS of Equation (5.31) leads to the minimax rate $n \exp(-(1 - \eta)\bar{n}_{\min}I)$. Up to a constant not dependent on π, Z^* or A , the

second term can be written as

$$\sum_{i=1}^n f_{i,2}(\pi - Z^*, A_{i,\cdot}) \lesssim \sum_a (\pi_{\cdot,a} - Z_{\cdot,a}^*)^T (A - \mathbb{E}A)(A - \mathbb{E}A)^T (\pi_{\cdot,a} - Z_{\cdot,a}^*).$$

In this way it is all about the random matrix $A - \mathbb{E}A$ and there exist sharp bounds on $\|A - \mathbb{E}A\|_{\text{op}}$. Note that $\sum_a \|\pi_{\cdot,a} - Z_{\cdot,a}^*\|^2 \leq \sum_a \|\pi_{\cdot,a} - Z_{\cdot,a}^*\|_1 \leq \|\pi - Z^*\|_1$. The second term ends up being upper bounded by $\|\pi - \pi^*\|_1$ multiplied by a coefficient factor.

Proof of Part One. Denote $z = r^{-1}(Z^*)$. By the definition of $h_{t,\lambda}(\cdot)$ in Equation (5.11), we have

$$\begin{aligned} \|[h_{t,\lambda}(\pi)]_{i,\cdot} - Z_{i,\cdot}^*\|_1 &\leq \frac{2 \sum_{a \neq z_i} \pi_{i,a}^{\text{pri}} \exp \left[2t \sum_{j \neq i} \pi_{j,a} (A_{i,j} - \lambda) \right]}{\sum_a \pi_{i,a}^{\text{pri}} \exp \left[2t \sum_{j \neq i} \pi_{j,a} (A_{i,j} - \lambda) \right]} \\ &\leq 2w \sum_{a \neq z_i} 1 \wedge \exp \left[2t \sum_{j \neq i} (\pi_{j,a} - \pi_{j,z_i}) (A_{i,j} - \lambda) \right]. \end{aligned}$$

Define $f(x) = 1 \wedge \exp(-x)$. It can be shown that for any $x_0 < 0$ and any integer $m \geq 1$ we have $f(x) \leq \exp(x_0) + \sum_{l=0}^{m-1} \exp(lx_0/m) \mathbb{I}\{x \geq (l+1)x_0/m\}$, which can be seen as a stepwise approximation of the continuous function $f(x)$. By taking $x_0 = -(n_a + n_{z_i})I/2$ and letting $x = 2t \sum_{j \neq i} (\pi_{j,a} - \pi_{j,z_i}) (A_{i,j} - \lambda)$, we have

$$\begin{aligned} \|[h_{t,\lambda}(\pi)]_{i,\cdot} - Z_{i,\cdot}^*\|_1 &\leq 2w \sum_{a \neq z_i} \exp \left[-\frac{(n_a + n_{z_i})I}{2} \right] + 2w \sum_{l=0}^{m-1} \left[\exp \left[-\frac{l(n_a + n_{z_i})I}{2m} \right] \right. \\ &\quad \left. \times \sum_{a \neq z_i} \mathbb{I} \left[2t \sum_{j \neq i} (\pi_{j,a} - \pi_{j,z_i}) (A_{i,j} - \lambda) \geq -\frac{(l+1)(n_a + n_{z_i})I}{2m} \right] \right]. \end{aligned}$$

We choose some $m \rightarrow \infty$ slowly such that

$$m = o(\bar{n}_{\min} I) \text{ and } m = o([wnI/[k[n/\bar{n}_{\min}]^2]^{1/4}]). \quad (5.32)$$

Thus, we have

$$\begin{aligned} \|h_{t,\lambda}(\pi) - Z^*\|_1 &\leq 2wnk \exp(-\bar{n}_{\min}I) + 2w \sum_{l=0}^{m-1} \sum_{a=1}^k \sum_{b \neq a} \left[\exp \left[-\frac{l(n_a + n_b)I}{2m} \right] \right. \\ &\quad \left. \times \sum_{i:z_i=b} \mathbb{I} \left[\sum_{j \neq i} (\pi_{j,a} - \pi_{j,b})(A_{i,j} - \lambda) \geq -\frac{(l+1)(n_a + n_b)I}{4mt} \right] \right] \end{aligned} \quad (5.33)$$

where we use the fact that $\min_{a \neq b} (n_a + n_b)/2 \geq \bar{n}_{\min}$.

The key to the rest of the analysis is to understand Equation (5.33) through the decomposition of the critical quantity $\sum_{j \neq i} (\pi_{j,a} - \pi_{j,b})(A_{i,j} - \lambda)$. We will show for any pair of $a, b \in [k]$ such that $a \neq b$, and any $i \in [n]$ such that $z_i = b$, it is equal to a summation of two terms: one only involves the ground truth Z^* , and the other involves the deviation $\pi - Z^*$. The former remains steady along iterations and contributes to the minimax rate, while the latter needs to be connected with the error $\|\pi - Z^*\|_1$.

Let $\theta_{a,b}$ be a vector of length n such that $[\theta_{a,b}]_j = \pi_{j,a} - Z_{j,a}^* + Z_{j,b}^* - \pi_{j,b}, \forall j \in [n]$. Then we have

$$\begin{aligned} \sum_{j \neq i} (\pi_{j,a} - \pi_{j,b})(A_{i,j} - \lambda) &= \sum_{j \neq i} (Z_{j,a}^* - Z_{j,b}^*)(A_{i,j} - \lambda) + \sum_{j \neq i} (\pi_{j,a} - Z_{j,a}^* + Z_{j,b}^* - \pi_{j,b})(A_{i,j} - \lambda) \\ &= \sum_{j \neq i} (Z_{j,a}^* - Z_{j,b}^*)(A_{i,j} - \lambda) + \sum_{j \neq i} (A_{i,j} - \lambda)[\theta_{a,b}]_j \\ &= \underbrace{\sum_{j \neq i} (Z_{j,a}^* - Z_{j,b}^*)(A_{i,j} - \lambda)}_{\text{involves } Z^*} + \underbrace{(A_{i,\cdot} - \mathbb{E}A_{i,\cdot})\theta_{a,b} + \sum_{j \neq i} (\mathbb{E}A_{i,j} - \lambda)[\theta_{a,b}]_j}_{\text{involves } \pi - Z^*}. \end{aligned} \quad (5.34)$$

With the help of Equation (5.34), Equation (5.33) can be written as

$$\begin{aligned}
& \|h_{t,\lambda}(\pi) - Z^*\|_1 \\
& \leq 2wnk \exp(-\bar{n}_{\min}I) + 2w \sum_{l=0}^{m-1} \sum_{a=1}^k \sum_{b \neq a} \left[\exp \left[-\frac{l(n_a + n_b)I}{2m} \right] \right. \\
& \quad \times \sum_{i:z_i=b} \mathbb{I} \left[\sum_{j \neq i} (Z_{j,a}^* - Z_{j,b}^*)(A_{i,j} - \lambda) \geq -\frac{(l+3/2)(n_a + n_b)I}{4mt} - \sum_{j \neq i} (\mathbb{E}A_{i,j} - \lambda)[\theta_{a,b}]_j \right] \\
& \quad \left. + 2w \sum_{a=1}^k \sum_{b \neq a} \left[\left[\sum_{l=0}^{m-1} \exp \left[-\frac{l(n_a + n_b)I}{2m} \right] \right] \times \sum_{i:z_i=b} \mathbb{I} \left[(A_{i,\cdot} - \mathbb{E}A_{i,\cdot})\theta_{a,b} \geq \frac{\bar{n}_{\min}I}{4mt} \right] \right].
\end{aligned}$$

Equations (5.18) and (5.32) imply $\sum_{l=0}^{m-1} \exp[-l(n_a + n_b)I/(2m)] \leq 2$. Thus, we have

$$\|h_{t,\lambda}(\pi) - Z^*\|_1 \leq 2wnk \exp(-\bar{n}_{\min}I) + \underbrace{2wL_1^{\text{sum}}}_{\text{involves } Z^*} + \underbrace{4wL_2^{\text{sum}}}_{\text{involves } \pi - Z^*},$$

where

$$L_1^{\text{sum}} \triangleq \sum_{l=0}^{m-1} \sum_{a=1}^k \sum_{b \neq a} \exp \left[-\frac{l(n_a + n_b)I}{2m} \right] \sum_{i:z_i=b} L_{1,i}(a, b, l),$$

with $L_{1,i}(a, b, l) \triangleq \mathbb{I}[\sum_{j \neq i} (Z_{j,a}^* - Z_{j,b}^*)(A_{i,j} - \lambda) \geq -(l+3/2)(n_a + n_b)I/(4mt) - \sum_{j \neq i} (\mathbb{E}A_{i,j} - \lambda)[\theta_{a,b}]_j]$, and

$$L_2^{\text{sum}} \triangleq \sum_{a=1}^k \sum_{b \neq a} \sum_{i:z_i=b} \mathbb{I} \left[(A_{i,\cdot} - \mathbb{E}A_{i,\cdot})\theta_{a,b} \geq \frac{\bar{n}_{\min}I}{4mt} \right].$$

In this way we turn $\|h_{t,\lambda}(\pi) - Z^*\|_1$ into calculations on L_1^{sum} and L_2^{sum} , where the former only involves the ground truth Z^* and the latter only involves the deviation $\pi - Z^*$.

We can obtain upper bounds on L_1^{sum} and L_2^{sum} as follows. Their proofs are deferred to the end of this section.

- For L_1^{sum} , there exists a sequence $\eta'' = o(1)$ such that with probability at least $1 - \exp[-2(\bar{n}_{\min}I)^{\frac{1}{2}}]$, we have

$$L_1^{\text{sum}} \leq nmk \exp[-(1 - 2\eta'')\bar{n}_{\min}I]. \tag{5.35}$$

- For L_2^{sum} , there exist constants c and r such that with probability at least $1 - n^{-r} - \exp(-5np^*)$, we have

$$L_2^{\text{sum}} \leq \frac{cknp^* \|\pi - Z^*\|_1}{(\bar{n}_{\min}I/(mt^*))^2} + \frac{cn^2kp^* \exp(-5np^*)}{\bar{n}_{\min}I/(mt^*)}. \quad (5.36)$$

Thus, we have

$$\begin{aligned} \|h_{t,\lambda}(\pi) - Z^*\|_1 &\leq 2wnk \exp(-\bar{n}_{\min}I) + 2wnmk \exp[-(1 - 2\eta'')\bar{n}_{\min}I] \\ &\quad + \frac{4cknp^* \|\pi - Z^*\|_1}{(\bar{n}_{\min}I/(mt^*))^2} + \frac{4cwk n^2 p^* \exp(-5np^*)}{\bar{n}_{\min}I/(mt^*)}, \end{aligned}$$

with probability at least $1 - \exp[-2(\bar{n}_{\min}I)^{\frac{1}{2}}] - n^{-r} - \exp(-5np^*)$. By Propositions 4.1 and 5.1, we have $p^*t^{*2} \asymp I$. Then due to Equation (5.32), we have

$$\frac{wkn p^*}{(\bar{n}_{\min}I/(mt^*))^2} \asymp wm^2 \left[\frac{n}{\bar{n}_{\min}} \right]^2 \frac{k}{nI} = o \left[\frac{1}{\sqrt{nI/[wk[n/\bar{n}_{\min}]^2]}} \right],$$

and

$$\frac{wkn^2 p^* \exp(-5np^*)}{\bar{n}_{\min}I/(mt^*)} \asymp wmk \frac{\sqrt{np^*}}{\sqrt{nI}} \left[\frac{n}{\bar{n}_{\min}} \right] n \exp(-5np^*) \leq n \exp(-5\bar{n}_{\min}I).$$

Thus, with probability at least $1 - \exp[-(\bar{n}_{\min}I)^{\frac{1}{2}}] - n^{-r}$, there exists some $\eta = o(1)$, such that

$$\|h_{t,\lambda}(\pi) - Z^*\|_1 \leq n \exp(-(1 - \eta)\bar{n}_{\min}I) + \frac{\|\pi - Z^*\|_1}{\sqrt{nI/[wk[n/\bar{n}_{\min}]^2]}}.$$

The proof for *Part One* is complete. The very last thing remained to be obtained is upper bounds on L_1^{sum} and L_2^{sum} , i.e., Equations (5.35) and (5.36). Recall the definition of $\theta_{a,b}$. We have some properties on $\theta_{a,b}$ which will be useful in the analysis for L_1^{sum} and L_2^{sum} : $\|\theta_{a,b}\|_{\infty} \leq 2$ and

$$\|\theta_{a,b}\|_1 \leq \|\pi_{\cdot,a} - Z^*_{\cdot,a}\|_1 + \|\pi_{\cdot,b} - Z^*_{\cdot,b}\|_1 \leq \|\pi - Z^*\|_1 \leq \gamma \bar{n}_{\min}, \quad (5.37)$$

and

$$\sum_{a=1}^k \sum_{b \neq a} \|\theta_{a,b}\|_1 \leq 2k \sum_a \|\pi_{\cdot,a} - Z_{\cdot,a}^*\|_1 \leq 2k \|\pi - Z^*\|_1. \quad (5.38)$$

1. *Bounds on L_1^{sum} .* By applying Markov inequality, we have

$$\begin{aligned} & \mathbb{E}L_{1,i}(a, b, l) \\ &= \mathbb{P} \left[t^* \sum_{j \neq i} (Z_{j,a}^* - Z_{j,b}^*)(A_{i,j} - \lambda) \geq -\frac{t^*(l+3/2)(n_a+n_b)I}{4mt} - t^* \sum_{j \neq i} (\mathbb{E}A_{i,j} - \lambda)[\theta_{a,b}]_j \right] \\ &\leq \exp \left[\frac{t^*(l+3/2)(n_a+n_b)I}{4mt} + t^*(\mathbb{E}A_{i,j} - \lambda 1_n^T) \theta_{a,b} \right] \mathbb{E} \exp \left[t^* \sum_{j \neq i} (Z_{j,a}^* - Z_{j,b}^*)(A_{i,j} - \lambda) \right]. \end{aligned}$$

With the help of Proposition 3.4, we have

$$\begin{aligned} & \mathbb{E} \exp \left[t^* \sum_{j \neq i} (Z_{j,a}^* - Z_{j,b}^*)(A_{i,j} - \lambda) \right] \\ &= \exp(-t^*(\lambda - \lambda^*)(n_a - n_b)) \exp(-t^*\lambda^*(n_a - n_b)) \prod_{j \neq i} \mathbb{E} \exp(t^*(Z_{j,a}^* - Z_{j,b}^*)A_{i,j}) \\ &= \exp(-t^*(\lambda - \lambda^*)(n_a - n_b)) \left[e^{-t\lambda} \frac{\mathbb{E}e^{tX}}{\mathbb{E}e^{-tY}} \right]^{\frac{n_a-n_b}{2}} \left[\mathbb{E}e^{tX} \mathbb{E}e^{-tY} \right]^{\frac{n_a+n_b}{2}} \\ &= \exp(-t^*(\lambda - \lambda^*)(n_a - n_b)) \exp \left[-\frac{(n_a + n_b)I}{2} \right]. \end{aligned}$$

Hence

$$\begin{aligned} & \mathbb{E}L_1^{\text{sum}} \quad (5.39) \\ &= \sum_{l=0}^{m-1} \sum_{a=1}^k \sum_{b \neq a} \left[\exp \left[-\frac{l(n_a+n_b)I}{2m} \right] \exp \left[\frac{t^*(l+3/2)(n_a+n_b)I}{4mt} + t^* \sum_{j \neq i} (\mathbb{E}A_{i,j} - \lambda)[\theta_{a,b}]_j \right] \right. \\ &\quad \left. \times \exp(-t^*(\lambda - \lambda^*)(n_a - n_b)) \exp \left[-\frac{(n_a+n_b)I}{2} \right] \right] \\ &\leq \sum_{l=0}^{m-1} \sum_{a=1}^k \sum_{b \neq a} \exp \left[-\frac{(1 + \frac{l}{m} - \frac{t^*(l+3/2)}{2mt})(n_a+n_b)I}{2} - t^*(\lambda - \lambda^*)(n_a - n_b) + t^* \sum_{j \neq i} (\mathbb{E}A_{i,j} - \lambda)[\theta_{a,b}]_j \right]. \end{aligned}$$

We are going to show $-(1 - \eta'')\bar{n}_{\min}I$ upper bounds terms in the exponent of RHS of Equation (5.39) by some $\eta'' = o(1)$. We first present some properties of λ^*, t^* and I that will be helpful:

$$I \asymp (p^* - q^*)^2/p^*, \quad (5.40)$$

$$\lambda^* \in (q^*, p^*), \quad (5.41)$$

$$\text{and } t^* \asymp (p^* - q^*)/p^*. \quad (5.42)$$

Here Equations (5.40) and (5.41) are proved by Propositions 4.1 and 5.1 respectively. Equation (5.42) is due to $t^* \asymp \log(1 + (p^* - q^*)/q^*) \asymp (p^* - q^*)/p^*$ under the assumption that $p^*, q^* = o(1)$, $p^* \asymp q^*$.

The first term in the exponent of Equation (5.39) is upper bounded by $-(1 - 7/(8m))\bar{n}_{\min}I$ by the assumption $t^*/t = 1 + o(1)$. Since $|t^*(\lambda - \lambda^*)| \leq \eta' t^*(p^* - q^*)$, by Equations (5.40) and (5.42) the second term is upper bounded by $\eta' \bar{n}_{\min}I$ up to a constant factor. For the last term in the exponent of Equation (5.39), since $|\lambda - \lambda^*| \leq \eta'(p^* - q^*)$ we have

$$\begin{aligned} t^* \left| \sum_{j \neq i} (\mathbb{E}A_{i,j} - \lambda)[\theta_{a,b}]_i \right| &\leq t^* \left| \sum_{j \neq i} (\mathbb{E}A_{i,j} - \lambda^*)[\theta_{a,b}]_i \right| + t^* \left| \sum_{j \neq i} (\lambda^* - \lambda)[\theta_{a,b}]_i \right| \\ &\leq (1 + \eta') t^*(p^* - q^*) \|\theta_{a,b}\|_1 \\ &\leq (1 + \eta') t^*(p^* - q^*) \gamma \bar{n}_{\min} \\ &\lesssim \gamma \bar{n}_{\min} I, \end{aligned}$$

where we use Equations (5.37) and (5.40) - (5.42).

As a consequence, there exists a sequence $\eta'' = o(1)$ that goes to zero slower than m^{-1}, γ, η' , such that the summation of three terms in the exponent of the RHS of Equation (5.39) is upper bounded by $-(1 - \eta'')\bar{n}_{\min}I$. Thus, Equation (5.39) can be written as

$$\mathbb{E}L_1^{\text{sum}} \leq nmk \exp[-(1 - \eta'')\bar{n}_{\min}I].$$

Since η'' goes to 0 slower than m^{-1} , we have $\eta'' \geq m^{-1} \geq (\bar{n}_{\min}I)^{\frac{1}{4}}$ by Equation (5.32).

Then by applying Markov inequality, we have

$$\mathbb{P} [L_1^{\text{sum}} \geq nmk \exp [-(1 - 2\eta'')\bar{n}_{\min}I]] \leq \exp [-\eta''\bar{n}_{\min}I] \leq \exp \left[-2(\bar{n}_{\min}I)^{\frac{1}{2}} \right].$$

That is, with probability at least $1 - \exp[-2(\bar{n}_{\min}I)^{\frac{1}{2}}]$, Equation (5.35) holds.

2. *Bounds on L_2^{sum} .* Depending on whether the network is dense or sparse, we consider two scenarios.

(1) *Dense Scenario:* $q^* \geq (\log n)/n$. In this scenario, we have a sharp bound on $\|A - \mathbb{E}A\|_{\text{op}}$. First we observe that

$$\begin{aligned} \sum_{i:z_i=b} [(A_{i,\cdot} - \mathbb{E}A_{i,\cdot})\theta_{a,b}]^2 &= \theta_{a,b}^T \sum_{i:z_i=b} [(A_{i,\cdot} - \mathbb{E}A_{i,\cdot})^T (A_{i,\cdot} - \mathbb{E}A_{i,\cdot})]\theta_{a,b} \\ &\leq \theta_{a,b}^T \sum_i [(A_{i,\cdot} - \mathbb{E}A_{i,\cdot})^T (A_{i,\cdot} - \mathbb{E}A_{i,\cdot})]\theta_{a,b} \\ &= \theta_{a,b}^T [(A - \mathbb{E}A)^T (A - \mathbb{E}A)]\theta_{a,b}. \end{aligned}$$

By applying Markov inequality, we have

$$L_2^{\text{sum}} \leq \sum_{a=1}^k \sum_{b \neq a} \frac{\theta_{a,b}^T [(A - \mathbb{E}A)^T (A - \mathbb{E}A)]\theta_{a,b}}{(\bar{n}_{\min}I/(4mt))^2}.$$

Since $\|\theta_{a,b}\|_{\infty} \leq 2$, we have $\|\theta_{a,b}\|^2 \leq 2\|\theta_{a,b}\|_1$. Lemma 5.4 shows $\|A - \mathbb{E}A\|_{\text{op}} \leq \sqrt{c_1 n p}$ holds with probability at least $1 - n^{-r}$ for some constants $c_1, r > 0$. Together with Equation (5.38), we have

$$\begin{aligned} \sum_{a=1}^k \sum_{b \neq a} \theta_{a,b}^T [(A - \mathbb{E}A)^T (A - \mathbb{E}A)]\theta_{a,b} &\leq \sum_{a=1}^k \sum_{b \neq a} \|A - \mathbb{E}A\|_{\text{op}}^2 \|\theta_{a,b}\|^2 \\ &\leq \sum_{a=1}^k \sum_{b \neq a} 2c_1 n p \|\theta_{a,b}\|_1 \\ &\leq 4c_1 k n p \|\pi - Z^*\|_1. \end{aligned}$$

Thus, with probability at least $1 - n^{-r}$,

$$L_2^{\text{sum}} \leq \frac{4c_1 knp \|\pi - Z^*\|_1}{(\bar{n}_{\min} I / (4mt))^2}.$$

(2) *Sparse Scenario:* $q^* < (\log n)/n$. When the network is sparse, the previous upper bound on $\|A - \mathbb{E}A\|_{\text{op}}$ no longer holds. Instead, removing nodes with large degrees is required to yield provably sharp bound on $\|A - \mathbb{E}A\|_{\text{op}}$. Define $S = \{i \in [n], \sum_j A_{i,j} \geq 20np^*\}$. We define \tilde{A}, \tilde{P} such that $\tilde{A}_{i,j} = A_{i,j} \mathbb{I}\{i, j \notin S\}$ and $\tilde{P}_{i,j} = (\mathbb{E}A_{i,j}) \mathbb{I}\{i, j \notin S\}$. Then we have the decomposition as

$$\begin{aligned} L_2(a, b) &\triangleq \sum_{i:z_i=b} \mathbb{I} \left[(A_{i,\cdot} - \mathbb{E}A_{i,\cdot}) \theta_{a,b} \geq \frac{\bar{n}_{\min} I}{4mt} \right] \\ &\leq \sum_{i:z_i=b} \mathbb{I} \left[(\tilde{A}_{i,\cdot} - \tilde{P}_{i,\cdot}) \theta_{a,b} \geq \frac{\bar{n}_{\min} I}{8mt} \right] \\ &\quad + \sum_{i:z_i=b} \mathbb{I} \left[\sum_{j \neq i} (A_{i,j} - \mathbb{E}A_{i,j}) [\theta_{a,b}]_{i,j} \mathbb{I}\{i \in S \text{ or } j \in S\} \geq \frac{\bar{n}_{\min} I}{8mt} \right] \\ &\triangleq L_{2,1}(a, b) + L_{2,2}(a, b). \end{aligned}$$

Define $L_{2,1}^{\text{sum}} \triangleq \sum_{a=1}^k \sum_{b \neq a} L_{2,1}(a, b)$. We have

$$L_{2,1}^{\text{sum}} \leq \sum_{a=1}^k \sum_{b \neq a} \frac{\theta_{a,b}^T [(\tilde{A} - \tilde{P})^T (\tilde{A} - \tilde{P})] \theta_{a,b}}{(\bar{n}_{\min} I / (8mt))^2} \leq \sum_{a=1}^k \sum_{b \neq a} \frac{2 \|\tilde{A} - \tilde{P}\|_{\text{op}}^2 \|\theta_{a,b}\|_1}{(\bar{n}_{\min} I / (8mt))^2}.$$

Lemma 4.2 shows $\|\tilde{A} - \tilde{P}\|_{\text{op}} \leq \sqrt{c_2 n p}$ holds with probability at least $1 - n^{-1}$ for some constant $c_2 > 0$. Then we have

$$L_{2,1}^{\text{sum}} \leq \frac{4c_2 knp \|\pi - Z^*\|_1}{(\bar{n}_{\min} I / (8mt))^2}.$$

Lemma 5.3 shows $\sum_{i,j} |A_{i,j} - \mathbb{E}A_{i,j}| \mathbb{I}\{i \in S\} \leq 20n^2 p^* \exp(-5np^*)$ holds with probability

at least $1 - \exp(-5np^*)$. Then by applying Markov inequality, we have

$$\begin{aligned}
L_{2,2}^{\text{sum}} &\triangleq \sum_{a=1}^k \left[\sum_{b \neq a} L_{2,2}(a,b) \right] \\
&\leq \sum_{a=1}^k \sum_{i,j=1}^n \frac{|A_{i,j} - \mathbb{E}A_{i,j}| |\theta_{a,b}]_{i,j}| \mathbb{I}\{i \in S \text{ or } j \in S\}}{\bar{n}_{\min} I / (8mt)} \\
&\leq \sum_{a=1}^k \frac{4 \sum_{i,j} |A_{i,j} - \mathbb{E}A_{i,j}| \mathbb{I}\{i \in S\}}{\bar{n}_{\min} I / (8mt)} \\
&\leq \frac{80n^2 k p^* \exp(-5np^*)}{\bar{n}_{\min} I / (8mt)}.
\end{aligned}$$

As a consequence, we have

$$L_2^{\text{sum}} \leq L_{2,1}^{\text{sum}} + L_{2,2}^{\text{sum}} \leq \frac{4c_2 k n p^* \|\pi - Z^*\|_1}{(\bar{n}_{\min} I / (8mt))^2} + \frac{80n^2 k p^* \exp(-5np^*)}{\bar{n}_{\min} I / (8mt)},$$

with probability at least $1 - n^{-1} - \exp(-5np^*)$. By the bounds on L_1^{sum} and L_2^{sum} , and due to $t/t^* = 1 + o(1)$, we obtain Equation (5.36).

5.5.4 Additional Lemmas and Propositions and Their Proofs

Lemma 5.1. *Let c_{init} be some sufficiently small constant. Consider any $\pi \in \Pi_1$ such that $\|\pi - Z^*\|_1 \leq c_{\text{init}} n/k$. Let $\alpha_p, \beta_p, \alpha_q, \beta_q, t, \lambda$ be the outputs after one step CAVI iteration from π described in Algorithm 4. That is, they are defined as Equations (5.26) - (5.29).*

Define

$$\hat{p} = \frac{\sum_{i < j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} A_{i,j}}{\sum_{i < j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a}}, \text{ and } \hat{q} = \frac{\sum_{i < j} \sum_{a \neq b} \pi_{i,a} \pi_{j,b} A_{i,j}}{\sum_{i < j} \sum_{a \neq b} \pi_{i,a} \pi_{j,b}}.$$

Under the same assumption as in Theorem 5.3, there exists some sequence $\epsilon = o(1)$ such that with probability at least $1 - e^3 5^{-n}$, the following inequality holds

$$\max \left\{ \frac{|\hat{p} - p^*|}{p^* - q^*}, \frac{|\hat{q} - q^*|}{p^* - q^*}, \frac{|t - t^*|}{(p^* - q^*)/p^*}, \frac{|\lambda - \lambda^*|}{p^* - q^*} \right\} \leq \epsilon + 24c_0 \frac{\|\pi - Z^*\|_1}{n/k},$$

uniformly over all the eligible π . In addition if we further assume c_{init} goes to 0, the LHS of the above inequality will be simply upper bounded by ϵ .

Proof. We are going to obtain tight bounds on $|\hat{p} - p^*|$ and $|\hat{q} - q^*|$ first. Note that we have the “variance-bias” decomposition as in

$$|\hat{p} - p^*| \leq \frac{|\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} (A_{i,j} - \mathbb{E}A_{i,j})|}{\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a}} + \left| \frac{\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} \mathbb{E}A_{i,j}}{\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a}} - p^* \right|.$$

We have concentration inequality holds for the numerator in the first term by Lemma 5.2.

That is, with probability at least $1 - e^{35-n}$, we have

$$\left| \sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} (A_{i,j} - \mathbb{E}A_{i,j}) \right| = \left| \frac{1}{2} \langle A - \mathbb{E}A, \pi \pi^T \rangle \right| \leq 3n \sqrt{np^*}$$

holds uniformly over all $\pi \in \Pi_1$. For the denominator, we have

$$\frac{n^2}{2} \geq \sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} = \frac{1}{2} \sum_{a=1}^k \|\pi_{\cdot,a}\|_1^2 \geq \frac{n^2}{2k},$$

since $\sum_{a=1}^k \|\pi_{\cdot,a}\|_1 = n$. Thus, we are able to obtain an upper bound on the first term as

$$\frac{|\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} (A_{i,j} - \mathbb{E}A_{i,j})|}{\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a}} \leq 6 \sqrt{\frac{k^2 p^*}{n}}.$$

For the second term, since $\mathbb{E}A_{i,j} = p^* \sum_{a=1}^k Z_{i,a}^* Z_{j,a}^* + q^* (1 - \sum_{a=1}^k Z_{i,a}^* Z_{j,a}^*)$, we have

$$\begin{aligned} \left| \frac{\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} \mathbb{E}A_{i,j}}{\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a}} - p^* \right| &= (p^* - q^*) \frac{\left| \sum_{i<j} \left[\sum_{a=1}^k \pi_{i,a} \pi_{j,a} \right] \left[\sum_{a=1}^k 1 - Z_{i,a}^* Z_{j,a}^* \right] \right|}{\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a}} \\ &= (p^* - q^*) \frac{|\langle \pi \pi^T, 11^T - Z^* Z^{*T} \rangle|}{\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a}} \\ &= (p^* - q^*) \frac{|\langle \pi \pi^T - Z^* Z^{*T}, 11^T - Z^* Z^{*T} \rangle|}{\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a}}, \end{aligned}$$

where in the last inequality we use the orthogonality between $Z^* Z^{*T}$ and $11^T - Z^* Z^{*T}$. For

its numerator, we have

$$\begin{aligned}
|\langle \pi \pi^T - Z^* Z^{*T}, 11^T - Z^* Z^{*T} \rangle| &\leq \|\pi \pi^T - Z^* Z^{*T}\|_1 \\
&\leq \|\pi - Z^*\|_1 (\|\pi\|_1 + \|Z^*\|_1) \\
&\leq \|\pi - Z^*\|_1 (2\|Z^*\|_1 + \|\pi - Z^*\|_1) \\
&\leq 3n \|\pi - Z^*\|_1.
\end{aligned}$$

This leads to

$$\left| \frac{\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} \mathbb{E} A_{i,j}}{\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a}} - p^* \right| \leq \frac{3n \|\pi - Z^*\|_1 (p^* - q^*)}{n^2/k} \leq 3kn^{-1} (p^* - q^*) \|\pi - Z^*\|_1.$$

Thus,

$$|\hat{p} - p^*| \leq 6\sqrt{\frac{k^2 p^*}{n}} + 3kn^{-1} (p^* - q^*) \|\pi - Z^*\|_1 \leq \left[\sqrt{\frac{k^2 p^*}{n(p^* - q^*)^2}} + \frac{3\|\pi - Z^*\|_1}{n/k} \right] (p^* - q^*).$$

Similar result holds for $|\hat{q} - q^*|$. Denote $\eta_0 = \sqrt{\frac{k^2 p^*}{n(p^* - q^*)^2}} + \frac{3\|\pi - Z^*\|_1}{n/k}$, thus

$$\max\{|\hat{p} - p^*|, |\hat{q} - q^*|\} \leq \eta_0 (p^* - q^*).$$

By the assumption of nI in Equation (5.18) and Proposition 4.1, we have $n(p^* - q^*)^2 / (k^2 p^*) \asymp nI/k^2 \rightarrow \infty$. Therefore, the first term in η_0 goes to 0. The second term in η_0 is at most $3c_{\text{init}}$ which implies $\eta_0 \leq 4c_{\text{init}}$.

By the fact that the digamma function satisfies $\psi(x) \in (\log(x - 1/2), \log x), \forall x \geq 1/2$, we have

$$\begin{aligned}
\psi(\alpha_p) - \psi(\beta_p) &\geq \log \frac{\alpha_p - 1/2}{\beta_p} \\
&= \log \left[\frac{\left[\alpha_p^{\text{pri}} - 1/2 + \sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} A_{i,j} \right] / \left[\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} \right]}{1 + \left[\beta_p^{\text{pri}} - \sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} A_{i,j} \right] / \left[\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} \right]} \right] \\
&= \log \left[\frac{\hat{p} + (\alpha_p^{\text{pri}} - 1/2) / \left[\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} \right]}{1 - \hat{p} + \beta_p^{\text{pri}} / \left[\sum_{i<j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} \right]} \right].
\end{aligned}$$

Recall that we have shown $\sum_{i < j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a}$ lies in the interval of $(n^2/(2k), n^2/2)$. By Equation (5.18), there exists a sequence $\eta' = o(1)$ such that $\alpha_p, \beta_p \leq \eta'(p^* - q^*)n^2/k$. Then we have

$$\psi(\alpha_p) - \psi(\beta_p) \geq \log \frac{p^* - |p^* - \hat{p}| - \eta'(p^* - q^*)}{1 - p^* + |p^* - \hat{p}| + \eta'(p^* - q^*)}.$$

Similar analysis leads to

$$\psi(\alpha_q) - \psi(\beta_q) \leq \log \frac{q^* + |q^* - \hat{q}| + \eta'(p^* - q^*)}{1 - q^* - |q^* - \hat{q}| - \eta'(p^* - q^*)}.$$

Together we have

$$\begin{aligned} t - t^* &\geq \log \left[\frac{p^* - |p^* - \hat{p}| - \eta'(p^* - q^*)}{1 - p^* + |p^* - \hat{p}| + \eta'(p^* - q^*)} \frac{1 - q^* - |q^* - \hat{q}| - \eta'(p^* - q^*)}{q^* + |q^* - \hat{q}| + \eta'(p^* - q^*)} \right] - t^* \\ &\geq \log \left[\left[1 - \frac{|p^* - \hat{p}| + \eta'(p^* - q^*)}{q^*} \right]^4 \frac{p^*(1 - q^*)}{q^*(1 - p^*)} \right] - t^* \\ &= 4 \log \left[1 - (\eta_0 + \eta') \frac{p^* - q^*}{q^*} \right]. \end{aligned}$$

Recall that we assume $c_0 p^* < q^* < p^*$. Thus $(\eta_0 + \eta')(p^* - q^*)/p^* \leq 5c_{\text{init}}c_0$. When c_{init} is sufficiently small, we have $(\eta_0 + \eta')(p^* - q^*)/p^* \leq 1/2$. Then using the fact $-x \geq \log(1-x) \geq -2x, \forall x \in (0, 1/2)$. We have

$$t - t^* \geq -8(\eta_0 + \eta')(p^* - q^*)/q^*.$$

Analogously we can obtain the same upper bound on $\hat{t} - t^*$, and then

$$|t - t^*| \leq 8c_0(\eta_0 + \eta') \frac{p^* - q^*}{p^*}.$$

Identical analysis can be applied towards bounds on $|\hat{\lambda} - \lambda^*|$. Note that

$$\log \frac{\beta_p}{\alpha_p + \beta_p} = \log \left[\frac{1 - \hat{p} + \beta_p^{\text{pri}} / \left[\sum_{i < j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} \right]}{1 + (\alpha_p^{\text{pri}} + \beta_p^{\text{pri}}) / \left[\sum_{i < j} \sum_{a=1}^k \pi_{i,a} \pi_{j,a} \right]} \right],$$

similarly for α_q, β_q . Omitting the immediate steps, we end up with

$$|\lambda - \lambda^*| = |[\psi(\beta_q) - \psi(\alpha_q + \beta_q)] - [\psi(\beta_p) - \psi(\alpha_p + \beta_p)] - \lambda^*| \leq 8(\eta_0 + \eta')(p^* - q^*).$$

The proof is complete after we unify and rephrase all the aforementioned results. \blacksquare

Lemma 5.2. *Let $A \in [0, 1]^{n \times n}$ such that $A = A^T$ and $A_{i,i} = 0, \forall i \in [n]$. Assume $\{A_{i,j}\}_{i < j}$ are independent random variable, and there exists $p \leq 1$ such that $9n^{-1} \leq \frac{2}{n(n-1)} \sum_{i < j} \text{Var}(A_{i,j}) \leq p$, and then we have*

$$\sup_{\pi \in \Pi_1} \left| \langle A - \mathbb{E}A, \pi \pi^T \rangle \right| \leq 6n\sqrt{np},$$

with probability at least $1 - e^{35^{-n}}$.

Proof. This result is a direct consequence of Grothendieck inequality [22] (see also Theorem 3.1 of [23] for a rephrased statement) on the matrix $A - \mathbb{E}A$. The Lemma 4.1 of [23] proves that with probability at least $1 - e^{35^{-n}}$,

$$\sup_{s,t \in \{-1,1\}^n} \left| \sum_{i,j} (A_{i,j} - \mathbb{E}A_{i,j}) s_i t_j \right| \leq 3n\sqrt{np}.$$

Then by applying Grothendieck inequality we obtain

$$\sup_{\|X_i\|_2 \leq 1, \forall i \in [n]} \left| \sum_{i,j} (A_{i,j} - \mathbb{E}A_{i,j}) X_i^T X_j \right| \leq 3cn\sqrt{np},$$

where c is a positive constant smaller than 2. This concludes with

$$\sup_{\pi \in \Pi_1} \left| \langle A - \mathbb{E}A, \pi \pi^T \rangle \right| \leq 6n\sqrt{np},$$

\blacksquare

Lemma 5.3. Let $A \in \{0, 1\}^{n \times n}$ be a symmetric binary matrix with $A_{i,i} = 0, \forall i \in [n]$, and $\{A_{i,j}\}_{i < j}$ are independent Bernoulli random variable. Let $p \geq \max_{i,j} \mathbb{E}A_{i,j}$. Define $S = \{i \in [n], \sum_j A_{i,j} \geq 20np\}$ and $Z_i = \sum_j |A_{i,j} - \mathbb{E}A_{i,j}| \mathbb{I}\{i \in S\}$. Then with probability at least $1 - \exp(-5np)$, we have

$$\sum_i Z_i \leq 20n^2p \exp(-5np).$$

Proof. Note that $\mathbb{E} \sum_j |A_{i,j} - \mathbb{E}A_{i,j}| \leq 2np(1-p) \leq 2np$. For any $s \geq 20np$, we have

$$\begin{aligned} \mathbb{P}(Z_i > s) &\leq \mathbb{P} \left[\sum_j |A_{i,j} - \mathbb{E}A_{i,j}| - \mathbb{E} \sum_j |A_{i,j} - \mathbb{E}A_{i,j}| > s - 2np \right] \\ &\leq \exp \left[-\frac{\frac{1}{2}(s - 2np)^2}{np + \frac{1}{3}(s - 2np)} \right] \\ &\leq \exp(-s/2), \end{aligned}$$

by implementing Bernstein inequality. Applying Bernstein inequality again we have

$$\begin{aligned} \mathbb{P}(Z_i > 0) &= \mathbb{P} \left[\sum_j A_{i,j} \geq 20np \right] \\ &\leq \mathbb{P} \left[\sum_j A_{i,j} - \mathbb{E} \sum_j A_{i,j} \geq 18np \right] \\ &\leq \exp \left[-\frac{(18np)^2/2}{np + 18np/3} \right] \\ &\leq \exp(-21np/2). \end{aligned}$$

Thus, we are able to bound $\mathbb{E}Z_i$ with

$$\begin{aligned} \mathbb{E}Z_i &\leq \int_0^{20np} \mathbb{P}(Z_i > 0) ds + \int_{20np}^{\infty} \mathbb{P}(Z_i > s) ds \\ &\leq 20np \exp(-21np/2) + \int_{20np}^{\infty} \exp(-s/2) ds \\ &\leq 20np \exp(-10np). \end{aligned}$$

By Markov inequality, we have

$$\begin{aligned} \mathbb{P} \left[\sum_{i,j} |A_{i,j} - \mathbb{E}A_{i,j}| \mathbb{I}\{i \in S\} \geq 20n^2 p \exp(-5np) \right] &= \mathbb{P} \left[\sum_i Z_i \geq 20n^2 p \exp(-5np) \right] \\ &\leq \frac{n\mathbb{E}Z_1}{20n^2 p \exp(-5np)} \\ &\leq \exp(-5np). \end{aligned}$$

■

Proposition 5.1. Define $\lambda = \log \frac{1-q}{1-p} / \log \frac{p(1-q)}{q(1-p)}$. For any $p, q > 0$ such that $p, q = o(1)$ and $p \asymp q$, there exists a constant $0 < c < 1/2$ such that

$$\frac{\lambda - q}{p - q} \in (c, 1 - c).$$

Proof. First we are going to establish the lower bound. Let $x = p - q$, and then we can rewrite λ as

$$\lambda = \frac{1}{1 + \frac{\log(1+x/q)}{\log(1+x/(1-q-x))}}.$$

Case I: $x \geq q/10$ Define $s = (p - q)/q$. Since $p \asymp q$ we have $s \geq 1/10$ and also upper bounded by some constant. We have

$$\begin{aligned} \frac{\lambda - q}{p - q} &= \frac{1}{s} \left[\frac{1}{q \left(1 + \frac{\log(1+s)}{\log(1+sq/(1-(s+1)q))} \right)} - 1 \right] \\ &= \frac{1}{s} \left[\frac{(1-q) \log(1+sq/(1-(s+1)q)) - q \log(1+s)}{q \log(1+sq/(1-(s+1)q)) + q \log(1+s)} \right] \\ &\geq \frac{1}{s} \frac{(1-q) \frac{sq}{1-(s+1)q} - q \log(1+s)}{2q \log(1+s)} \\ &\geq \frac{1}{8} \frac{1-q}{\log(1+s)}, \end{aligned}$$

which is lower bounded by some constant $c > 0$.

Case II: $x < q/10$ By Taylor theorem, there exist $0 \leq \epsilon_1, \epsilon_2 \leq 1/10$ such that

$$\log \left[1 + \frac{x}{q} \right] = \frac{x}{q} - \frac{1 - \epsilon_1}{2} \left[\frac{x}{q} \right]^2,$$

$$\text{and } \log \left[1 + \frac{x}{1 - q - x} \right] = \frac{x}{1 - q - x} - \frac{1 - \epsilon_2}{2} \left[\frac{x}{1 - q - x} \right]^2.$$

Thus, we have

$$\frac{\log(1 + \frac{x}{q})}{\log(1 + \frac{x}{1-q-x})} = \frac{q(1-q)^2 - [2q(1-q) + \frac{1-\epsilon_1}{2}(1-q)^2]x + c_1x^2 + c_2x^3}{q^2(1-q) - \frac{3-\epsilon_2}{2}q^2x},$$

where $c_1 = (1 - \epsilon_1)(1 - q) + q$ and $c_2 = -(1 - \epsilon_1)/2$. Thus,

$$\begin{aligned} \frac{\lambda - q}{p - q} &= \frac{1}{x} \left[\frac{q^2(1-q) - \frac{3-\epsilon_2}{2}q^2x}{q(1-q) - [2q(1-q) + \frac{1-\epsilon_1}{2}(1-q)^2 + \frac{3-\epsilon_2}{2}q^2]x + c_1x^2 + c_2x^3} - q \right] \\ &= \frac{[\frac{1}{2}q(1-q) + \frac{\epsilon_2}{2}q^2(1-q) - \frac{\epsilon_1}{2}(1-q)^2q] + c_1qx + c_2qx^2}{q(1-q) - [2q(1-q) + \frac{1-\epsilon_1}{2}(1-q)^2 + \frac{3-\epsilon_2}{2}q^2]x + c_1x^2 + c_2x^3} \end{aligned}$$

Note that $|c_1|, |c_2| \leq 1$. We have

$$\frac{\lambda - q}{p - q} \geq \frac{\frac{1}{4}q(1-q)}{2q(1-q)} \geq 1/8.$$

By using exactly the same discussion, we can show $(p - \lambda)/(p - q) > c$. Thus, we proved the desired bound stated in the proposition. ■

Lemma 5.4. [Theorem 5.2 of [35]] Let $A \in \{0, 1\}^{n \times n}$ be a symmetric binary matrix with $A_{i,i} = 0, \forall i \in [n]$, and $\{A_{i,j}\}_{i < j}$ are independent Bernoulli random variable. If $p \triangleq \max_{i,j} \mathbb{E}A_{i,j} \geq \log n/n$. Then there exist constants $c, r > 0$ such that

$$\|A - \mathbb{E}A\|_{\text{op}} \leq c\sqrt{np},$$

with probability at least $1 - n^{-r}$.

Chapter 6

Generalization: Degree-corrected Block Model

Despite a rich literature dedicated to their theoretical properties, SBMs suffer significant drawbacks when it comes to modeling real world social and biological networks. In particular, due to the model assumption, all nodes within the same community in an SBM are exchangeable and hence have the same degree distribution. In comparison, nodes in real world networks often exhibit degree heterogeneity even when they belong to the same community [43]. For example, Bickel and Chen [8] showed that for a karate club network, SBM does not provide a good fit for the data set, and the resulting clustering analysis is qualitatively different from the truth.

One way to accommodate degree heterogeneity is to introduce a set of degree-correction parameters $\{\theta_i : i = 1, \dots, n\}$, one for each node, which can be interpreted as the popularity or importance of a node in the network. Then one could revise the edge distributions to $A_{ij} = A_{ji} \stackrel{ind.}{\sim} \text{Bern}(\theta_i \theta_j B_{z(i)z(j)})$ for all $i > j$, and this gives rise to the Degree-Corrected Block Model (DCBM) [14, 33]. In a DCBM, within the same community, a node with a larger value of degree-correction parameter is expected to have more connections than that with a smaller value. On the other hand, SBMs are special cases of DCBM in which the degree-correction parameters are all equal. Empirically, the larger class of DCBM is able to provide possibly much better fits to many real world network datasets [43]. Throughout

the chapter, we allow k and B to scale with n as n tends to infinity. Since the proposal of the model, there have been various methods proposed for community detection in DCBM, including but not limited to spectral clustering [25, 30, 35, 44] and modularity based approaches [5, 11, 33, 54]. On the theoretical side, [24] provides an information-theoretic characterization of the impossibility region of community detection for DCBM with two clusters, and sufficient conditions have been given in [11, 54] for strongly and weakly consistent community detection. However, two fundamental statistical questions remain unanswered:

- What are the fundamental limits of community detection in DCBM?
- Once we know these limits, can we achieve them adaptively via some polynomial time algorithm?

These two questions are also the main topics of Chapter 3 and 4 for the SBM. We are going to provide answers for them in this chapter for the DCBM.

6.1 Model

Recall that a random graph of size n generated by a DCBM has its adjacency matrix A satisfying $A_{ii} = 0$ for all $i \in [n]$ and

$$A_{ij} = A_{ji} \stackrel{ind}{\sim} \text{Bern}(\theta_i \theta_j B_{z(i)z(j)}) \text{ for all } i \neq j \in [n]. \quad (6.1)$$

For each $u \in [k]$ and a given $z \in [k]^n$, we let $n_u = n_u(z) = \sum_{i=1}^n \mathbf{1}_{\{z(i)=u\}}$ be the size of the u^{th} community. Let $P = \mathbb{E}[A] \in [0, 1]^{n \times n}$. We propose to consider the following parameter

space for DCBM of size n :

$$\begin{aligned}
\mathcal{P}_n(\theta, p, q, k, \beta; \delta) = & \{P \in [0, 1]^{n \times n} : \exists z \in [k]^n \text{ and } B = B^T \in \mathbb{R}^{k \times k}, \\
& \text{s.t. } P_{ii} = 0, P_{ij} = \theta_i \theta_j B_{z(i)z(j)}, \forall i \neq j \in [n], \\
& \frac{1}{n_u} \sum_{z(i)=u} \theta_i \in [1 - \delta, 1 + \delta], \forall u \in [k], \\
& \max_{u \neq v} B_{uv} \leq q < p \leq \min_u B_{uu}, \\
& \frac{n}{\beta k} - 1 \leq n_u \leq \frac{\beta n}{k} + 1, \forall u \in [k]\}.
\end{aligned} \tag{6.2}$$

We are mostly interested in the behavior of minimax risks over a sequence of such parameter spaces as n tends to infinity and the key model parameters θ, p, q, k scale with n in some appropriate way. On the other hand, we take $\beta \geq 1$ as an absolute constant and require the (slack) parameter δ to be an $o(1)$ sequence throughout the chapter.

To see the rationale behind the definition in (6.2), let us examine each of the parameters used in the definition. The starting point is $\theta \in \mathbb{R}_+^n$, which we treat for now as a given sequence of degree-correction parameters. Given θ , we consider all possible label vectors z such that the approximate normalization $\frac{1}{n_u} \sum_{z(i)=u} \theta_i = 1 + o(1)$ holds for all communities. The introduction of the slack parameter $0 < \delta = o(1)$ rules out those parameter spaces in which community detection can be trivially achieved by only examining the normalization of the θ_i 's. On the other hand, the proposed normalization ensures that for all $u \neq v \in [k]$,

$$B_{uu} \approx \frac{1}{n_u(n_u - 1)} \sum_{i:z(i)=u} \sum_{j \neq i:z(j)=u} P_{ij} \quad \text{and} \quad B_{uv} \approx \frac{1}{n_u n_v} \sum_{i:z(i)=u} \sum_{j:z(j)=v} P_{ij}.$$

Therefore, B_{uu} and B_{uv} can be understood as the (approximate) average connectivity within the u^{th} community and between the u^{th} and the v^{th} communities, respectively. Under this interpretation, p can be seen as a lower bound on the within community connectivities and q an upper bound on the between community connectivities. We require the assumption $p > q$ to ensure that the model is ‘‘assortative’’ in an average sense. Finally, we also require the individual community sizes to be contained in the interval $[n/(\beta k) - 1, \beta n/k + 1]$. In other words, the community sizes are assumed to be of the same order. Although we have

focused on the case of assortative networks, we expect the same expression of minimax rates to hold in the disassortative case, i.e., $\min_{u \neq v} B_{uv} \geq q > p \geq \max_u B_{uu}$.

6.2 Minimax Risks

The key information-theoretic quantity that governs the minimax risk of community detection is I , which is defined through

$$\exp(-I) = \begin{cases} \frac{1}{n} \sum_{i=1}^n \exp(-\theta_i \frac{n}{2} (\sqrt{p} - \sqrt{q})^2), & k = 2, \\ \frac{1}{n} \sum_{i=1}^n \exp\left(-\theta_i \frac{n}{\beta^k} (\sqrt{p} - \sqrt{q})^2\right), & k \geq 3. \end{cases} \quad (6.3)$$

Note that it is very similar to Equation (3.2), thus we use the same notation I despite a bit abuse of notation.

Minimax upper bounds Given any parameter space $\mathcal{P}_n(\theta, p, q, k, \beta; \delta)$, we can define the following estimator:

$$\hat{z} = \underset{z' \in \mathcal{P}_n(\theta, p, q, k, \beta; \delta)}{\operatorname{argmax}} \prod_{1 \leq i < j \leq n} [(\theta_i \theta_j p)^{A_{ij}} (1 - \theta_i \theta_j p)^{1 - A_{ij}} \mathbf{1}_{\{z'(i) = z'(j)\}} + (\theta_i \theta_j q)^{A_{ij}} (1 - \theta_i \theta_j q)^{1 - A_{ij}} \mathbf{1}_{\{z'(i) \neq z'(j)\}}]. \quad (6.4)$$

If there is a tie, we break it arbitrarily. The estimator (6.4) is the maximum likelihood estimator for a special case of DCBM where $B_{uu} = p$ and $B_{uv} = q$ for all $u \neq v \in [k]$. In other cases, the objective function in (6.4) is a misspecified likelihood function. For any sequences $\{a_n\}$ and $\{b_n\}$, we write $a_n = \Omega(b_n)$ if $a_n \geq C b_n$ for some absolute constant $C > 0$ for all $n \geq 1$. The following theorem characterizes the asymptotic behavior of the risk bounds for the estimator (6.4).

Theorem 6.1 (Minimax Upper Bounds). *Consider any sequence*

$\{\mathcal{P}_n(\theta, p, q, k, \beta; \delta)\}_{n=1}^{\infty}$ *such that as* $n \rightarrow \infty$, $I \rightarrow \infty$, $p > q$, $\|\theta\|_{\infty} = o(n/k)$, $\min_{i \in [n]} \theta_i = \Omega(1)$ *and* $\log k = o(\min(I, \log n))$. *When* $k \geq 3$, *further assume* $\beta \in [1, \sqrt{5/3}]$. *Then the*

estimator in (6.4) satisfies

$$\limsup_{n \rightarrow \infty} \frac{1}{I} \log \left(\sup_{\mathcal{P}_n(\theta, p, q, k, \beta; \delta)} \mathbb{E} \ell(\hat{z}, z) \right) \leq -1.$$

Before proceeding, we briefly discuss the conditions in Theorem 6.1. First, the condition $\min_{i \in [n]} \theta_i = \Omega(1)$ requires that all θ_i 's are at least of constant order. One should note that this condition does not rule out the possibility that $\max_i \theta_i \gg \min_i \theta_i$, and so a great extent of degree variation, even within the same community, is allowed. Next, $\log k = o(\log n)$ requires that the number of communities k , if it diverges to infinity, grows at a sub-polynomial rate with the number of nodes n . Furthermore, $\beta \in [1, \sqrt{5/3})$ is a technical condition that we need for a combinatorial argument in the proof to go through when $k \geq 3$. When $k = O(1)$ and $\Omega(1) = \min_i \theta_i \leq \|\theta\|_\infty = O(1)$, Theorem 6.1 only requires $I \rightarrow \infty$, which is equivalent to $n(p - q)^2/p \rightarrow \infty$. Informed readers might find the result in Theorem 6.1 in parallel to that in [52]. However, due to the presence of degree-correction parameters, the proof of Theorem 6.1 is significantly different from that of the corresponding result in [52]. For example, a new folding argument is employed to deal with degree heterogeneity.

Minimax lower bounds We now show that the rates in Theorem 6.1 are asymptotic minimax optimal by establishing matching minimax lower bounds. To this end, we require the following condition on the degree-correction parameters $\theta \in \mathbb{R}_+^n$. The condition guarantees that $\mathcal{P}_n(\theta, p, q, k, \beta; \delta)$ is non-empty. Moreover, it is only needed for establishing minimax lower bounds.

Condition 1. We say that $\theta \in \mathbb{R}_+^n$ satisfies Condition N if

1. When $k = 2$, there exists a disjoint partition $\mathcal{C}_1, \mathcal{C}_2$ of $[n]$, such that $|\mathcal{C}_1| = \lfloor n/2 \rfloor$, $|\mathcal{C}_2| \in \{\lfloor n/2 \rfloor, \lfloor n/2 \rfloor + 1\}$ and $|\mathcal{C}_u|^{-1} \sum_{i \in \mathcal{C}_u} \theta_i \in (1 - \delta/4, 1 + \delta/4)$ for $u = 1, 2$.
2. When $k \geq 3$, there exists a disjoint partition $\{\mathcal{C}_u\}_{u \in [k]}$ of $[n]$, such that $|\mathcal{C}_1| \leq |\mathcal{C}_2| \leq \dots \leq |\mathcal{C}_k|$, $|\mathcal{C}_1| = |\mathcal{C}_2| = \lfloor n/(\beta k) \rfloor$ and $|\mathcal{C}_u|^{-1} \sum_{i \in \mathcal{C}_u} \theta_i \in (1 - \delta/4, 1 + \delta/4)$ for all $u \in [k]$.

We note that the condition is only on θ (as opposed to the parameter space) and the actually communities in the data generating model need not coincide with the partition

that occurs in the statement of the condition.

With the foregoing definition, we have the following result.

Theorem 6.2 (Minimax Lower Bounds). *Consider any sequence*

$\{\mathcal{P}_n(\theta, p, q, k, \beta; \delta)\}_{n=1}^\infty$ *such that as* $n \rightarrow \infty$, $I \rightarrow \infty$, $1 < p/q = O(1)$, $p\|\theta\|_\infty^2 = o(1)$, $\log k = o(I)$, $\log(1/\delta) = o(I)$ *and* θ *satisfies Condition N. Then*

$$\liminf_{n \rightarrow \infty} \frac{1}{I} \log \left(\inf_{\hat{z}} \sup_{\mathcal{P}_n(\theta, p, q, k, \beta; \delta)} \mathbb{E} \ell(\hat{z}, z) \right) \geq -1.$$

Compared with the conditions in Theorem 6.1, the conditions of Theorem 6.2 are slightly different. The condition $1 < p/q = O(1)$ ensures that the smallest average within community connectivity is of the same order as (albeit larger than) the largest average between community connectivity. Such an assumption is typical in the statistical literature on block models. The condition $\|\theta\|_\infty^2 p = o(1)$ ensures that the maximum expected node degree scales at a sublinear rate with the network size n . Furthermore, when $k = O(1)$, the condition $\log k = o(I)$ can be dropped because it is equivalent to $I \rightarrow \infty$, which in turn is necessary for the minimax risk to converge to zero.

Combining both theorems, we have the minimax risk of the problem.

Theorem 6.3. *Under the conditions of Theorems 6.1 and 6.2, we have*

$$\inf_{\hat{z}} \sup_{\mathcal{P}_n(\theta, p, q, k, \beta; \delta)} \mathbb{E} \ell(\hat{z}, z) = \exp(-(1 + o(1))I),$$

where $o(1)$ stands for a sequence whose absolute values tend to zero as n tends to infinity.

When $\theta = 1_n$, the foregoing minimax risk reduces to the corresponding result for SBM in the sparse regime where $q < p = o(1)$. In this case, Equation (6.3) implies that the minimax risk is

$$\exp(-(1 + o(1))I) = \begin{cases} \exp\left(- (1 + o(1)) \frac{n}{2} (\sqrt{p} - \sqrt{q})^2\right), & k = 2, \\ \exp\left(- (1 + o(1)) \frac{\beta n}{k} (\sqrt{p} - \sqrt{q})^2\right), & k \geq 3. \end{cases}$$

Note that when $q < p = o(1)$, the Rényi divergence of order $\frac{1}{2}$ used in the minimax risk

expression in Theorem 3.1 is equal to $(1+o(1))(\sqrt{p}-\sqrt{q})^2$, indicating a match with Theorem 3.1.

6.3 An Adaptive and Computationally Feasible Procedure

Theorem 6.1 shows that the minimax rate can be achieved by the estimator (6.4) obtained via combinatorial optimization which is not computationally feasible. Moreover, the procedure depends on the knowledge of the parameters θ , p and q . These features make it not applicable in practical situations. In this section, we introduce a two-stage algorithm for community detection in DCBM which is not only computationally feasible but also adaptive over a wide range of unknown parameter values. We show that the procedure achieves minimax optimal rates under certain regularity conditions.

6.3.1 A Two-Stage Algorithm

We first give the method for initialization (Algorithm 5), and then present the complete algorithm (Algorithm 6).

Initialization: weighted k -medians clustering. We first give Algorithm 5, which is an analogous of the low-rank based spectral clustering (i.e., Algorithm 1) for the regular SBM.

To explain the rationale behind our proposal, with slight abuse of notation, let $P = (P_{ij}) \in [0, 1]^{n \times n}$, where for all $i, j \in [n]$, $P_{ij} = P_{ji} = \theta_i \theta_j B_{z(i)z(j)}$. Except for the diagonal entries, P is the same as in (6.2). For any $i \in [n]$, let P_i denote the i^{th} row of P . Then for all i such that $z(i) = u$, we observe that

$$\theta_i^{-1} P_i = (\theta_1 B_{u,z(1)}, \dots, \theta_n B_{u,z(n)})$$

are all equal. Thus, there are exactly k different vectors that the normalized row vectors $\{\theta_i^{-1} P_i\}_{i=1}^n$ can be. Moreover, which one of the k vectors the i^{th} normalized row vector equals is determined solely by its community label $z(i)$. This observation suggests one can

Algorithm 5: Weighted k -medians Clustering

Data: Adjacency matrix $A \in \{0, 1\}^{n \times n}$, number of clusters k , tuning parameter τ .

Result: Initial label estimator \hat{z}^0 .

- 1 Define $T_\tau(A) \in \{0, 1\}^{n \times n}$ by replacing the i th row and column of A whose row sum is larger than τ by zeroes for each $i \in [n]$;
- 2 Solve

$$\hat{P} = \operatorname{argmin}_{\operatorname{rank}(P) \leq k} \|T_\tau(A) - P\|_{\mathbb{F}}^2;$$

- 3 Let \hat{P}_i be the i^{th} row of \hat{P} . Define $S_0 = \{i \in [n] : \|\hat{P}_i\|_1 = 0\}$. Set $\hat{z}^0(i) = 0$ for $i \in S_0$, and define $\tilde{P}_i = \hat{P}_i / \|\hat{P}_i\|_1$ for $i \notin S_0$;
- 4 Solve a $(1 + \epsilon)$ - k -median optimization problem on S_0^c . That is, find $\{\hat{z}^0(i)\}_{i \in S_0^c}$ in $[k]^{|S_0^c|}$ that satisfies

$$\sum_{u=1}^k \min_{v_u \in \mathbb{R}^n} \sum_{\{i \in S_0^c : \hat{z}^0(i) = u\}} \|\hat{P}_i\|_1 \|\tilde{P}_i - v_u\|_1 \leq (1 + \epsilon) \min_{z \in [k]^n} \sum_{u=1}^k \min_{v_u \in \mathbb{R}^n} \sum_{\{i \in S_0^c : z(i) = u\}} \|\hat{P}_i\|_1 \|\tilde{P}_i - v_u\|_1. \quad (6.5)$$

design a reasonable community detection procedure by clustering the sample counterparts of the vectors $\{\theta_1^{-1}P_1, \theta_2^{-1}P_2, \dots, \theta_n^{-1}P_n\}$, which leads us to the proposal of Algorithm 5.

In Algorithm 5, Steps 1 and 2 aim to find an estimator \hat{P} of P by solving a low rank approximation problem. Then, in Step 3, we can use $\|\hat{P}_i\|_1^{-1} \hat{P}_i$ as a surrogate for $\theta_i^{-1}P_i$. Finally, Step 4 performs a weighted k -median clustering procedure applied on the row vectors of the $n \times k$ matrix.

Full algorithm. The full algorithm for community detection in DCBM is given in Algorithm 6. It is analogous to Algorithm 2 for the regular SBM, with the only difference that we use normalized majority voting in Algorithm 6 instead of penalized majority voting as in Algorithm 2.

6.3.2 Performance Guarantees

In this part, we state high probability performance guarantees for the proposed procedure. The theoretical property of the algorithms requires an extra bound on the maximal entry

Algorithm 6: A Two-stage Algorithm for DCBM

- Data:** Adjacency matrix $A \in \{0, 1\}^{n \times n}$ and number of clusters k ;
Result: Clustering label estimator $\hat{z} \in [k]^n$;
- 1 For each $i \in [n]$, apply Algorithm 5 to A_{-i} . The result, which is a vector of dimension $n - 1$, is stored in $(\hat{z}_{-i}^0(1), \dots, \hat{z}_{-i}^0(i - 1), \hat{z}_{-i}^0(i + 1), \dots, \hat{z}_{-i}^0(n))$;
 - 2 For each $i \in [n]$, the i th entry of \hat{z}_{-i}^0 is set as

$$\hat{z}_{-i}^0(i) = \operatorname{argmax}_{u \in [k]} \frac{1}{|\{j : \hat{z}_{-i}^0(j) = u\}|} \sum_{j: \hat{z}_{-i}^0(j) = u} A_{ij};$$

- 3 Set $\hat{z}(1) = \hat{z}_{-1}^0(1)$. For each $i \in \{2, \dots, n\}$, set

$$\hat{z}(i) = \operatorname{argmax}_{u \in [k]} |\{j : \hat{z}_{-1}^0(j) = u\} \cap \{j : \hat{z}_{-i}^0(j) = \hat{z}_{-i}^0(i)\}|. \quad (6.6)$$

of EA. We incorporate this condition into the following parameter space

$$\begin{aligned} & \mathcal{P}'_n(\theta, p, q, k, \beta; \delta, \alpha) \\ &= \{P = (\theta_i \theta_j B_{z(i)z(j)} \mathbf{1}_{\{i \neq j\}}) \in \mathcal{P}_n(\theta, p, q, k, \beta; \delta) : \max_{u \in [k]} B_{uu} \leq \alpha p\}. \end{aligned}$$

The parameter α is assumed to be a constant no smaller than 1 that does not change with n . By studying the proofs of Theorem 6.2 and Theorem 6.1, the minimax lower and upper bounds do not change for the slightly smaller parameter space $\mathcal{P}'_n(\theta, p, q, k, \beta; \delta, \alpha)$. Therefore, the rate $\exp(-(1 + o(1))I)$ still serves as a benchmark for us to develop theoretically justifiable algorithms for the parameter space $\mathcal{P}'_n(\theta, p, q, k, \beta; \delta, \alpha)$.

Error rate for the initialization stage. As a first step, we provide the following high probability error bound for Algorithm 5.

Theorem 6.4 (Error Bound for Algorithm 5). *Assume $\delta = o(1)$, $1 < p/q = O(1)$ and $\|\theta\|_\infty = o(n/k)$. Let $\tau = C_1(np\|\theta\|_\infty^2 + 1)$ for some sufficiently large constant $C_1 > 0$ in Algorithm 5. Then, there exist some constants $C', C > 0$, such that for any generative model in $\mathcal{P}'_n(\theta, p, q, k, \beta; \delta, \alpha)$, we have with probability at least $1 - n^{-(1+C')}$,*

$$\min_{\rho} \sum_{\{i: \hat{z}(i) \neq \rho(z(i))\}} \theta_i \leq C \frac{(1 + \epsilon)k^{5/2} \sqrt{n\|\theta\|_\infty^2 p + 1}}{p - q},$$

where the minimization is over all the permutations on $[k]$.

Theorem 6.4 provides a uniform high probability bound for the sum of θ_i 's of the nodes which are assigned wrong labels. In the special case when $\theta_i = 1, \forall i \in [1]$ (i.e., under the regular SBM), it immediately implies $\ell(\hat{z}, z) \lesssim (1 + \epsilon)k^{5/2}\sqrt{np+1}/(p-q)$, a slightly weaker result compared to Theorem 4.1. The proof of Theorem 6.4 essentially follows that of Theorem 4.1, with additional effort to handle θ . Thus we omit it in this thesis and refer the readers to our paper [20] for details.

Error rate for the refinement stage. We now state a general high probability error bound for Algorithm 6. To introduce this result, we define another information-theoretic quantity. For any $t \in (0, 1)$, define

$$J_t(p, q) = 2(tp + (1-t)q - p^t q^{1-t}). \quad (6.7)$$

By Jensen's inequality, it is straightforward to verify that $J_t(p, q) \geq 0$ and $J_t(p, q) = 0$ if and only if $p = q$. As a special case, when $t = \frac{1}{2}$, we have

$$J_{\frac{1}{2}}(p, q) = (\sqrt{p} - \sqrt{q})^2. \quad (6.8)$$

For a given $z \in [k]^n$, let $n_{(1)} \leq \dots \leq n_{(k)}$ be the order statistics of community sizes $\{n_u(z) : u = 1, \dots, k\}$. Then, we define the quantity J by through

$$\exp(-J) = \frac{1}{n} \sum_{i=1}^n \exp\left(-\theta_i \left(\frac{n_{(1)} + n_{(2)}}{2}\right) J_{t^*}(p, q)\right) \quad (6.9)$$

with $t^* = \frac{n_{(1)}}{n_{(1)} + n_{(2)}}$. With the foregoing definitions, the following theorem gives a general error bound for Algorithm 6.

Theorem 6.5. *Under the conditions of Theorem 6.4, we further assume that $\delta = o(\frac{p-q}{p})$,*

$$\|\theta\|_\infty^2 p \geq n^{-1},$$

$$\frac{(1 + \epsilon)k^{5/2} \|\theta\|_\infty \sqrt{p}}{\sqrt{n}(p - q)} = o\left(\frac{p - q}{kp}\right), \quad \text{and} \quad (6.10)$$

$$\min_{\gamma \geq 0} \left\{ n^{-1} |\{i \in [n] : \theta_i \leq \gamma\}| + \frac{(1 + \epsilon)k^{5/2} \|\theta\|_\infty \sqrt{p}}{\gamma \sqrt{n}(p - q)} \right\} = o\left(\frac{p - q}{k^2 p}\right). \quad (6.11)$$

Then there is a sequence $\eta = o(1)$ such that the output \hat{z} of Algorithm 6 satisfies

$$\lim_{n \rightarrow \infty} \inf_{\mathcal{P}'_n(\theta, p, q, k, \beta; \delta, \alpha)} \mathbb{P} \{ \ell(\hat{z}, z) \leq \exp(- (1 - \eta)J) \} = 1.$$

Theorem 6.5 gives a general error bound for the performance of Algorithm 6. It shows that Algorithm 6 converges at the rate $\exp(-(1 + o(1))J)$. According to the properties of $J_t(p, q)$ stated in Appendix B of [20], one can show that when $n_{(1)} = (1 + o(1))n_{(2)}$, $J = (1 + o(1))I$, and that in general

$$n_{(1)}(\sqrt{p} - \sqrt{q})^2 \leq \left(\frac{n_{(1)} + n_{(2)}}{2}\right) J_{t^*}(p, q) \leq \left(\frac{n_{(1)} + n_{(2)}}{2}\right) (\sqrt{p} - \sqrt{q})^2.$$

Using this relation, we can state the convergence rate in Theorem 6.5 using the quantity I .

Corollary 6.1. *Under the conditions of Theorem 6.5, there is a sequence $\eta = o(1)$ such that the output \hat{z} of Algorithm 6 satisfies*

$$\begin{aligned} \lim_{n \rightarrow \infty} \inf_{\mathcal{P}'_n(\theta, p, q, 2, \beta; \delta, \alpha)} \mathbb{P} \{ \ell(\hat{z}, z) \leq \exp(- (1 - \eta)\beta^{-1}I) \} &= 1, \\ \lim_{n \rightarrow \infty} \inf_{\mathcal{P}'_n(\theta, p, q, k, \beta; \delta, \alpha)} \mathbb{P} \{ \ell(\hat{z}, z) \leq \exp(- (1 - \eta)I) \} &= 1, \text{ for } k \geq 3. \end{aligned}$$

Therefore, when $k \geq 3$, the minimax rate $\exp(-(1 + o(1))I)$ is achieved by Algorithm 6. The only situation where the minimax rate is not achieved by Algorithm 6 is when $k = 2$ and $\beta > 1$. For this case, there is an extra β^{-1} factor on the exponent of the convergence rate. The proof of Theorem 6.3 essentially follows that of Theorem 2 for the regular SBM. Hence we omitted it in this thesis and we refer readers to our paper [20] for details.

Bibliography

- [1] Emmanuel Abbe and Colin Sandon. Community detection in general stochastic block models: Fundamental limits and efficient algorithms for recovery. In *Foundations of Computer Science (FOCS), 2015 IEEE 56th Annual Symposium on*, pages 670–688. IEEE, 2015.
- [2] Lada A Adamic and Natalie Glance. The political blogosphere and the 2004 us election: divided they blog. In *Proceedings of the 3rd international workshop on Link discovery*, pages 36–43. ACM, 2005.
- [3] Edoardo M Airoldi, David M Blei, Stephen E Fienberg, and Eric P Xing. Mixed membership stochastic blockmodels. *Journal of Machine Learning Research*, 9(Sep):1981–2014, 2008.
- [4] Réka Albert, Hawoong Jeong, and Albert-László Barabási. Internet: Diameter of the world-wide web. *nature*, 401(6749):130, 1999.
- [5] Arash A Amini, Aiyou Chen, Peter J Bickel, and Elizaveta Levina. Pseudo-likelihood methods for community detection in large sparse networks. *The Annals of Statistics*, 41(4):2097–2122, 2013.
- [6] Albert-Laszlo Barabasi and Zoltan N Oltvai. Network biology: understanding the cell’s functional organization. *Nature Reviews Genetics*, 5(2):101–113, 2004.
- [7] Peter Bickel, David Choi, Xiangyu Chang, and Hai Zhang. Asymptotic normality of maximum likelihood and its variational approximation for stochastic blockmodels. *The Annals of Statistics*, 41(4):1922–1943, 2013.

- [8] Peter J Bickel and Aiyou Chen. A nonparametric view of network models and Newman–Girvan and other modularities. *Proceedings of the National Academy of Sciences*, 106(50):21068–21073, 2009.
- [9] Alain Celisse, Jean-Jacques Daudin, and Laurent Pierre. Consistency of maximum-likelihood and variational estimators in the stochastic block model. *Electronic Journal of Statistics*, 6:1847–1899, 2012.
- [10] Yudong Chen and Jiaming Xu. Statistical-computational tradeoffs in planted problems and submatrix localization with a growing number of clusters and submatrices. *The Journal of Machine Learning Research*, 17(1):882–938, 2016.
- [11] Yudong Chen, Xiaodong Li, and Jiaming Xu. Convexified modularity maximization for degree-corrected stochastic block models. *The Annals of Statistics (To Appear)*, 2018.
- [12] Peter Chin, Anup Rao, and Van Vu. Stochastic block model and community detection in sparse graphs: A spectral algorithm with optimal rate of recovery. In *Proceedings of The 28th Conference on Learning Theory*, pages 391–423, 2015.
- [13] Amin Coja-Oghlan. Graph partitioning via adaptive spectral techniques. *Combinatorics, Probability and Computing*, 19(02):227–284, 2010.
- [14] Anirban Dasgupta, John E Hopcroft, and Frank McSherry. Spectral analysis of random graphs with skewed degree distributions. In *Foundations of Computer Science, 2004. Proceedings. 45th Annual IEEE Symposium on*, pages 602–610. IEEE, 2004.
- [15] David Easley and Jon Kleinberg. *Networks, crowds, and markets: Reasoning about a highly connected world*. Cambridge University Press, 2010.
- [16] Yingjie Fei and Yudong Chen. Exponential error rates of SDP for block models: Beyond Grothendieck’s inequality. *arXiv preprint arXiv:1705.08391*, 2017.
- [17] Donniell E Fishkind, Daniel L Sussman, Minh Tang, Joshua T Vogelstein, and Carey E Priebe. Consistent adjacency-spectral partitioning for the stochastic block model when

- the model parameters are unknown. *SIAM Journal on Matrix Analysis and Applications*, 34(1):23–39, 2013.
- [18] Chao Gao, Aad W van der Vaart, and Harrison H Zhou. A general framework for bayes structured linear models. *arXiv preprint arXiv:1506.02174*, 2015.
- [19] Chao Gao, Zongming Ma, Anderson Y Zhang, and Harrison H Zhou. Achieving optimal misclassification proportion in stochastic block model. *The Journal of Machine Learning Research*, 18(60):1–45, 2017.
- [20] Chao Gao, Zongming Ma, Anderson Y Zhang, and Harrison H Zhou. Community detection in degree-corrected block models. *The Annals of Statistics (To Appear)*, 2018.
- [21] Michelle Girvan and Mark EJ Newman. Community structure in social and biological networks. *Proceedings of the National Academy of Sciences*, 99(12):7821–7826, 2002.
- [22] Alexandre Grothendieck. Résumé de la théorie métrique des produits tensoriels topologiques. *Resenhas do Instituto de Matemática e Estatística da Universidade de São Paulo*, 2(4):401–481, 1996.
- [23] Olivier Guédon and Roman Vershynin. Community detection in sparse networks via Grothendieck’s inequality. *Probability Theory and Related Fields*, 165(3-4):1025–1049, 2016.
- [24] Lennart Gulikers, Marc Lelarge, and Laurent Massoulié. An impossibility result for reconstruction in a degree-corrected planted-partition model. *arXiv preprint arXiv:1511.00546*, 2015.
- [25] Lennart Gulikers, Marc Lelarge, and Laurent Massoulié. A spectral method for community detection in moderately-sparse degree-corrected stochastic block models. *arXiv preprint arXiv:1506.08621*, 2015.
- [26] Bruce Hajek, Yihong Wu, and Jiaming Xu. Achieving exact cluster recovery threshold via semidefinite programming: Extensions. *IEEE Transactions on Information Theory*, 62(10):5918–5937, 2016.

- [27] Mark S Handcock, Adrian E Raftery, and Jeremy M Tantrum. Model-based clustering for social networks. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 170(2):301–354, 2007.
- [28] Jake M Hofman and Chris H Wiggins. Bayesian approach to network modularity. *Physical review letters*, 100(25):258701, 2008.
- [29] Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic block-models: First steps. *Social networks*, 5(2):109–137, 1983.
- [30] Jiashun Jin. Fast community detection by SCORE. *The Annals of Statistics*, 43(1):57–89, 2015.
- [31] Antony Joseph and Bin Yu. Impact of regularization on spectral clustering. *The Annals of Statistics*, 44(4):1765–1791, 2016.
- [32] Brian Karrer and Mark EJ Newman. Stochastic blockmodels and community structure in networks. *Physical Review E*, 83(1):016107, 2011.
- [33] Brian Karrer and Mark EJ Newman. Stochastic blockmodels and community structure in networks. *Physical Review E*, 83(1):016107, 2011.
- [34] Can M Le, Elizaveta Levina, and Roman Vershynin. Sparse random graphs: Regularization and concentration of the Laplacian. *arXiv preprint arXiv:1502.03049*, 2015.
- [35] Jing Lei and Alessandro Rinaldo. Consistency of spectral clustering in stochastic block models. *The Annals of Statistics*, 43(1):215–237, 2015.
- [36] László Lovász. *Large networks and graph limits*, volume 60. American Mathematical Soc., 2012.
- [37] Elchanan Mossel, Joe Neeman, and Allan Sly. Stochastic block models and reconstruction. *arXiv preprint arXiv:1202.1499*, 2012.
- [38] Elchanan Mossel, Joe Neeman, and Allan Sly. Consistency thresholds for the planted bisection model. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pages 69–75. ACM, 2015.

- [39] Elchanan Mossel, Joe Neeman, and Allan Sly. A proof of the block model threshold conjecture. *Combinatorica*, pages 1–44, 2017.
- [40] Mark Newman. *Networks: an introduction*. Oxford University Press, 2010.
- [41] Mark EJ Newman. The structure and function of complex networks. *SIAM review*, 45(2):167–256, 2003.
- [42] Mark EJ Newman and Elizabeth A Leicht. Mixture models and exploratory analysis in networks. *Proceedings of the National Academy of Sciences*, 104(23):9564–9569, 2007.
- [43] Tiago P Peixoto. Model selection and hypothesis testing for large-scale network models with overlapping groups. *Physical Review X*, 5(1):011033, 2015.
- [44] Tai Qin and Karl Rohe. Regularized spectral clustering under the degree-corrected stochastic blockmodel. In *Advances in Neural Information Processing Systems*, pages 3120–3128, 2013.
- [45] Zahra S Razaee, Arash A Amini, and Jingyi Jessica Li. Matched bipartite block model with covariates. *arXiv preprint arXiv:1703.04943*, 2017.
- [46] Karl Rohe, Sourav Chatterjee, and Bin Yu. Spectral clustering and the high-dimensional stochastic blockmodel. *The Annals of Statistics*, 39(4):1878–1915, 2011.
- [47] Daniel L Sussman, Minh Tang, Donniell E Fishkind, and Carey E Priebe. A consistent adjacency spectral embedding for stochastic blockmodel graphs. *Journal of the American Statistical Association*, 107(499):1119–1128, 2012.
- [48] Aad W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- [49] Piet Van Mieghem. *Performance analysis of communications networks and systems*. Cambridge University Press, 2006.
- [50] VAN VU. A simple svd algorithm for finding hidden partitions. *Combinatorics, Probability and Computing*, 27(1):124–140, 2018. doi: 10.1017/S0963548317000463.

- [51] Stanley Wasserman. *Social network analysis: Methods and applications*, volume 8. Cambridge university press, 1994.
- [52] Anderson Y Zhang and Harrison H Zhou. Minimax rates of community detection in stochastic block models. *The Annals of Statistics*, 44(5):2252–2280, 2016.
- [53] Anderson Y Zhang and Harrison H Zhou. Theoretical and computational guarantees of mean field variational inference for community detection. *arXiv preprint arXiv:1710.11268*, 2017.
- [54] Yunpeng Zhao, Elizaveta Levina, and Ji Zhu. Consistency of community detection in networks under degree-corrected stochastic block models. *The Annals of Statistics*, 40(4):2266–2292, 2012.