EXACT COMMUNITY RECOVERY UNDER SIDE INFOR MATION: OPTIMALITY OF SPECTRAL ALGORITHMS

Anonymous authors

004

006

008

009 010

011

012

013

014

015

016

017

018

019

021

024

025 026

028

Paper under double-blind review

ABSTRACT

We study the problem of exact community recovery in general, two-community block models, in the presence of node-attributed side information. We allow for a very general side information channel for node attributes, and for pairwise (edge) observations, consider both Bernoulli and Gaussian matrix models, capturing the Stochastic Block Model, Submatrix Localization, and \mathbb{Z}_2 -Synchronization as special cases. A recent work of Dreveton et al. (2024) characterized the informationtheoretic limit of a very general exact recovery problem with side information. In this paper, we show algorithmic achievability in the above important cases by designing a simple but optimal spectral algorithm that incorporates side information (when present) along with the eigenvectors of the pairwise observation matrix. Using the powerful tool of entrywise eigenvector analysis (Abbe et al., 2020), we show that our spectral algorithm can mimic the so called *genie-aided estimators*, where the i^{th} genie-aided estimator optimally computes the estimate of the i^{th} label, when all remaining labels are revealed by a genie. This perspective provides a unified understanding of the optimality of spectral algorithms for various exact recovery problems in a recent line of work.

027 1 INTRODUCTION

In this paper, we consider inference problems of the following form: there is an unknown partition of the set $[n] := \{1, 2, ..., n\}$ into two *communities*, denoted (C_+, C_-) with $(\rho, 1 - \rho)$ fraction of vertices respectively for $\rho \in (0, 1)$. The community assignments are encoded by a vector $\sigma^* \in \{\pm 1\}^n$. We observe a symmetric matrix $A \in \mathbb{R}^{n \times n}$, which is specified by three distributions: \mathcal{P}_+ , \mathcal{P}_- , and \mathcal{Q} . The entries of A are independent (up to symmetry), such that $A_{ij} \sim \mathcal{P}_+$ if $i, j \in C_+$, $A_{ij} \sim \mathcal{P}_-$ if $i, j \in C_-$, and $A_{ij} \sim \mathcal{Q}$ otherwise. One famous example is the Stochastic Block Model (SBM), where $\mathcal{P}_+ \equiv \text{Bern}(p_1)$, $\mathcal{P}_- \equiv \text{Bern}(p_2)$, and $\mathcal{Q} \equiv \text{Bern}(q)$. Other prominent examples include \mathbb{Z}_2 -synchronization and submatrix localization, in which $\mathcal{P}_+, \mathcal{P}_-, \mathcal{Q}$ are Gaussian distributions. Given the observation A, the goal is to *exactly* recover the unknown σ^* .

After a long line of research (Decelle et al., 2011; Mossel et al., 2015; Abbe et al., 2016; Hajek et al., 2016; Abbe, 2017; Bandeira et al., 2017; Javanmard et al., 2016; Cai et al., 2017), a fairly complete picture of the fundamental limits of exact recovery and algorithmic acheivability is known. In recent 040 years, with practical considerations, exact recovery has also been studied in different variants such as 041 node-attributed side information (Saad & Nosratinia, 2018; 2020; Deshpande et al., 2018; Dreveton 042 et al., 2024), partially censored edges (Abbe et al., 2014; Hajek et al., 2015; Moghaddam et al., 043 2022; Dhara et al., 2022a; 2023), multiple correlated networks (Racz & Sridhar, 2021; Gaudio et al., 2022), spatially embedded networks (Abbe et al., 2021; Gaudio et al., 2024b;a), etc. See Section 044 1.2 for different variants considered more broadly in community detection literature beyond exact recovery. 046

In this paper, similar to Dreveton et al. (2024), we consider the setting of node-attributed side information, where we are given a side information vector $y \in \mathcal{Y}^n$, in addition to the pairwise observation matrix A. The side information is drawn according to a pair of distributions $(\mathcal{S}_+, \mathcal{S}_-)$, where for each $i \in [n]$, we obtain

051

- $y_i \sim \mathcal{S}_+ \text{ if } i \in C_+, \quad \text{ or } \quad y_i \sim \mathcal{S}_- \text{ if } i \in C_-,$
- 053 where the observations are independent, conditioned on the communities. The setup captures several interesting special cases:

- Gaussain Features (GF): For each node, we have a vector of d real-valued attributes. This is modeled as 𝒴 = ℝ^d and both 𝒴₊ and 𝒴₋ are parameterized multivariate Gaussians in d dimensions.
- 2. Binary Erasure Channel (BEC): The true community assignment of some subset of vertices are known (partially *observed* labels), modeled as $\mathcal{Y} = \{-1, 0, +1\}$ and y is formed by passing σ^* through a binary erasure channel.
- 3. Binary Symmetric Channel (BSC): We have a "guess" on community assignment of each vertex (partially *correct* labels) and interested in recovering the true assignment. This is modeled as passing σ^* through a binary symmetric channel with $\mathcal{Y} \in \{-1, +1\}$.

Exact recovery under BEC and BSC side information channels was first studied by Saad & Nos-064 ratinia (2018; 2020) for two special cases of Bernoulli pairwise observations; symmetric SBM 065 $(\mathcal{P}_+ \equiv \mathcal{P}_-, \rho = 1/2)$ and Planted Dense Subgraph (PDS) $(\mathcal{P}_- \equiv \mathcal{Q})$. The work of Dreveton et al. 066 (2024) recently studied a very general recovery problem, allowing generic distributions ($\mathcal{P}_+, \mathcal{P}_-, \mathcal{Q}$) 067 and side information laws $(\mathcal{S}_+, \mathcal{S}_-)$, satisfying certain technical assumptions. They derived sharp 068 information theoretic thresholds for exact recovery under side information in a unified way, and 069 showed that the optimal Maximum A Posterior (MAP) estimator succeeds down to this threshold. However, naïve MAP estimation requires a brute-force computation, and thus it is not efficiently 071 (poly-time) computable in the worst case. Therefore, it remains important to design efficient algo-072 rithms. We note that Dreveton et al. also proposed an efficient, iterative likelihood maximization 073 algorithm when node and edge observations are from an exponential family of distributions, though without theoretical guarantees on the performance. In this work, our primary goal is: 074

075

055

056

058

059

060

061

062

063

076 077

085

087

Objective 1: To design efficient algorithms that are provably optimal i.e. succeed down to the information theoretic threshold.

We note that the work of Saad & Nosratinia (2018; 2020) designed provably optimal algorithms only for the two special cases they studied (symmetric SBM and PDS), and only under BEC and BSC channels. They use a two-stage strategy that first determines *almost exactly correct* labels, followed by a refinement step achieving exact recovery. In this work, our focus will be on designing a *singlestage* spectral algorithm that directly achieves exact recovery and in much greater generality. In particular, we consider any side information channel with distributions (S_+ , S_-) but restrict to the following pairwise distribution laws¹:

- 1. Rank One Spike (ROS): The distributions \mathcal{P}_+ , \mathcal{P}_- , and \mathcal{Q} are Gaussian distributions, capturing \mathbb{Z}_2 -synchronization and submatrix localization as special cases.
- 2. SBM: The general two community stochastic block model where \mathcal{P}_+ , \mathcal{P}_- , and \mathcal{Q} are any Bernoulli distributions.

Background on spectral algorithms for exact recovery. Spectral algorithms were popularized by classical works such as McSherry's algorithm for community detection (McSherry, 2001) and 091 the planted clique recovery algorithm of Alon, Krivilevich, and Sudakov (Alon et al., 1998). In 092 recent years, attention has shifted to spectral algorithms without the need of a combinatorial cleanup phase. This line of work was initiated by the influential work of Abbe, Fan, Wang, and Zhong (Abbe 094 et al., 2020), who showed that in the symmetric, balanced SBM, simply thresholding the second leading eigenvector u_2 of the adjacency matrix at 0 gives the correct community partition with high 096 probability. In order to analyze the vector u_2 , Abbe et al. developed the technique of *entrywise* 097 eigenvector analysis, which characterizes the entrywise behavior of eigenvectors of matrices whose 098 expectations are low-rank.

Following Abbe et al. (2020), a series of papers used the entrywise eigenvector technique to give strong guarantees for spectral algorithms. For example, Deng et al. (2021) showed that, in the symmetric SBM, using the Laplacian instead of the adjacency matrix also yields an optimal algorithm. Another line of work Dhara et al. (2022a; 2023) considered the censored variant of the problem, where the status of some edges is unknown. Even for the censored variant, they show that

¹⁰⁴ ¹While our side information channel laws are general, we note that the information theoretic limits of ¹⁰⁵ Dreveton et al. (2024) are in even more generality. Namely, they consider (i) more general distribution families ¹⁰⁶ ($\mathcal{P}_+, \mathcal{P}_-, \mathcal{Q}$), and (ii) the case of more than two communities. In Appendix A.1, we shall describe how our ¹⁰⁷ spectral algorithms can be further generalized to (i) any ($\mathcal{P}_+, \mathcal{P}_-, \mathcal{Q}$) coming from exponential families (with ¹⁰⁸ additional requirements), and (ii) more than two communities.

108 spectral algorithms are optimal, where the encoding of the unknown edges is chosen carefully to 109 achieve optimaliy. Perhaps surprisingly, Dhara et al. (2023) showed that for this censored variant 110 of the problem, to handle cases beyond the symmetric SBM ($\mathcal{P}_+ \equiv \mathcal{P}_-, \rho = 1/2$) and the PDS 111 $(\mathcal{P}_{-} \equiv \mathcal{Q})$, any clustering algorithm based on a single adjacency matrix does not succeed down to 112 the information-theoretic threshold. Instead, the authors devised a spectral algorithm which forms two matrix representations of the same network and takes a carefully chosen linear combination of 113 their eigenvectors to achieve optimaliy. All these results at their core rely on the entrywise behavior 114 of eigenvectors (Abbe et al., 2020). 115

This raises several closely related questions: what governs the optimality of these seemingly different problem specific choices? Is there some principle behind designing new algorithms in related settings? For example, for the standard uncensored variant as we study here, the spectral algorithm for general two community SBM is unknown even when there is no side information. Do we need two matrices like its censored counterpart (Dhara et al., 2023)? To answer these questions, another auxiliary goal of this work:

Objective 2: To develop a unified perspective on the optimality of these spectral algorithms.

123 124 1.1 Our Contribution

125 **Spectral Algorithms.** We propose a simple spectral strategy of computing the top eigenvectors 126 (top one for ROS and top two for SBM) of the observation matrix A and take an appropriate linear 127 combination. When side information y is also available, we incorporate it by shifing the eigenvector 128 combination by the log-likelihood ratio vector of side information:

$$\log\left(\frac{\mathcal{S}_{+}}{\mathcal{S}_{-}}(y)\right) \in (\mathbb{R} \cup \{\pm\infty\})^{n} \text{ such that } \left[\log\left(\frac{\mathcal{S}_{+}}{\mathcal{S}_{-}}(y)\right)\right]_{i} := \log\left(\frac{\mathcal{S}_{+}(y_{i})}{\mathcal{S}_{-}(y_{i})}\right), \tag{1}$$

followed by a prescribed thresholding. Our algorithm is highly efficient with nearly linear runtime for the respective setting.

The main technical novelty lies in establishing a rigorous connection between the spectral estimator 134 and the so-called *genie-aided estimators*, where i^{th} genie-aided estimator optimally computes the 135 $i^{\rm th}$ label, where the rest of them, denoted by σ_{-i}^* , are revealed by a genie. Using the entrywise 136 eigenvector technique (Abbe et al., 2020), we show that taking an appropriately weighted sum of 137 the leading eigenvectors of A (along with side information use prescribed above when they are 138 available) produces a vector whose i^{th} entry is well-approximated by the statistic computed by these 139 genie-aided estimators in each of the settings. Thus, the spectral algorithm without any clean-up step 140 is able to mimic the genie, and therefore achieves exact recovery down to the information-theoretic 141 limits. 142

Models without Side Information. We note that even in the case of no side information, determining when a single stage spectral algorithm without any clean up stage can succeed is of interest to the learning theory community. Our results fill the complete picture in important remaining cases such as the general SBM (beyond the symmetric case with $\mathcal{P}_+ \equiv \mathcal{P}_-$, $\rho = 1/2$ (Abbe et al., 2020) and the Planted Dense Subgraph (PDS) with $\mathcal{P}_- \equiv \mathcal{Q}$ (Dhara et al., 2022a)). Most notably, unlike its censored counterpart (Dhara et al., 2023), one does not need two matrix representations to achieve optimality. In fact, the strategies of Abbe et al. (2020); Deng et al. (2021); Dhara et al. (2022a;b; 2023) are essentially mimicking the genie in their respective settings.

150 151

152

122

129 130 131

1.2 RELATIONSHIP TO PRIOR WORK

Local to global amplification and two stage algorithms. The SBM is known to exhibit *local to-global amplification*, in the sense that whenever (local) recovery of a single vertex label given the labels of all other vertices is possible with probability 1 - o(1/n), then (global) exact recovery is possible with probability 1 - o(1) (see e.g., Abbe (2017)). Two-stage algorithms, which are prevalent in the literature (Abbe et al., 2016; Coja-Oghlan, 2006; Vu, 2018; Yun & Proutiere, 2014; 2016; Gaudio et al., 2024a), essentially leverage this statistical property.

158

Unsupervised vs Semi-supervised Learning. The presence of side information turns community recovery from an unsupervised learning problem into a semi-supervised learning problem. Therefore, another related perspective is to investigate the effectiveness semi-supervised learning approaches over unsupervised ones e.g. (Jiang & Ke, 2023; Wu et al., 2022; Ni et al., 2024). We

note that, from this perspective, investigating other relaxed goals such as *weak* or *almost exact* recovery could provide a more satisfying picture, as the exact recovery is very demanding criterion.
The primary message of our work is that, for exact recovery goal, there are simple extensions of the spectral algorithm that optimally combines the signal from A and the side information y, achieving the new sharp information-theoretic limit.

168 1.3 ORGANIZATION

Section 2 contains our models and other preliminary setup. Our main results are stated in Section
The genie-aided estimators are introduced in Section 4, which we show how to mimic using a
spectral strategy in Section 5. Future directions are proposed in Section 6. The proofs are postponed
to the appendices.

173 174

175

185

187 188

197 198

199

200 201 202

215

167

2 PRELIMINARIES

176 2.1 MODELS

We first introduce the General Two Community Block Model (GBM), which captures the two special cases that we consider.

Definition 2.1 (General Two Community Block Model (GBM)). For any $\rho \in (0, 1)$ and distributions $\mathcal{P}_+, \mathcal{P}_-$, and \mathcal{Q} , we say that $(A, \sigma^*) \sim \text{GBM}_n(\rho, \mathcal{P}_+, \mathcal{P}_-, \mathcal{Q})$, where $A \in \mathbb{R}^{n \times n}$ and $\sigma^* \in \{\pm 1\}^n$ are sampled as follows. Each coordinate of σ^* is sampled i.i.d. such that $\mathbb{P}(\sigma_i^* = +1) = \rho$ and $\mathbb{P}(\sigma_i^* = -1) = 1 - \rho$. Moreover, we will use the notation $C_+ := \{i : \sigma_i^* = +1\}$ and $C_- := \{i : \sigma_i^* = +1\}$. Conditioned on σ^* , we sample A, a zero diagonal symmetric matrix with independent entries, such that for $1 \le i < j \le n$, we have

$$A_{ij} \sim \begin{cases} \mathcal{P}_+, & \text{if } i, j \in C_+\\ \mathcal{P}_-, & \text{if } i, j \in C_-\\ \mathcal{Q}, & \text{otherwise.} \end{cases}$$

In the above definition, we restrict to distributions which are either (i) a continuous distribution or (ii) a finite, discrete distribution. In the above definition, we restrict ourselves to settings where distributions $\mathcal{P}_+(\cdot), \mathcal{P}_-(\cdot), \mathcal{Q}(\cdot)$ identify the corresponding probability density function or probability mass function, respectively. We will consider special cases where the distributions $(\mathcal{P}_+, \mathcal{P}_-, \mathcal{Q})$ are either all Gaussian or Bernoulli distributions. The specialized definitions are given below.

Definition 2.2 (Rank One Spike (ROS)). Fix any $\rho \in (0, 1)$ and $a, b \in \mathbb{R}$ such that $\max\{|a|, |b|\} > 0$. 195 0. We say that $(A, \sigma^*) \sim \text{ROS}_n(\rho, a, b)$ if they are sampled as follows. First sample σ^* as mentioned 196 for GBM. Conditioned on σ^* consider the vector $v^* \in \{a, b\}^n$ such that for $i \in [n]$

$$v_i^* = a \cdot \mathbf{1}[\sigma_i^* = +1] + b \cdot \mathbf{1}[\sigma_i^* = -1].$$
⁽²⁾

;;

Finally, conditioned on σ^* , we get independent noisy measurements for every i < j of the following form.

$$A_{ij} = v_i^* v_j^* \sqrt{\frac{\log n}{n}} + W_{ij}, \text{ where } W_{ij} \overset{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1).$$

Note that the model $\text{ROS}_n(\rho, a, b)$ is a special case of $\text{GBM}_n(\rho, \mathcal{P}_+, \mathcal{P}_-, \mathcal{Q})$, by taking

$$\mathcal{P}_{+} \equiv \mathcal{N}\left(a^{2}\sqrt{\frac{\log n}{n}}, 1\right), \mathcal{P}_{-} \equiv \mathcal{N}\left(b^{2}\sqrt{\frac{\log n}{n}}, 1\right), \text{ and } \mathcal{Q} \equiv \mathcal{N}\left(ab\sqrt{\frac{\log n}{n}}, 1\right).$$

Taking b = 0 yields a version of the Gaussian submatrix localization problem Hajek et al. (2018), for which the goal is to recover a submatrix of elevated mean (corresponding to the entries in C_+). Taking a = -b yields the \mathbb{Z}_2 -synchronization problem Bandeira et al. (2017) after rescaling; see Remark A.1. Our scaling choice allows both submatrix localization and \mathbb{Z}_2 -synchronization to be handled under a unified model. We also consider the Stochastic Block Model (SBM) in the logarithmic-degree regime, which is the relevant regime for exact recovery.

Definition 2.3 (Stochastic Block Model (SBM)). Fix any $\rho \in (0, 1)$ and $a_1, a_2, b > 0$. Then the model SBM_n(ρ, a_1, a_2, b) is a special case of GBM_n($\rho, \mathcal{P}_+, \mathcal{P}_-, \mathcal{Q}$) with

$$\mathcal{P}_{+} \equiv \operatorname{Bern}\left(\frac{a_{1}\log n}{n}\right), \mathcal{P}_{-} \equiv \operatorname{Bern}\left(\frac{a_{2}\log n}{n}\right), and \mathcal{Q} \equiv \operatorname{Bern}\left(\frac{b\log n}{n}\right)$$

Finally, we consider a generic side information channel.

Definition 2.4 (Side Information (SI)). For any domain \mathcal{Y} , distributions $(\mathcal{S}_+, \mathcal{S}_-)$ supported on \mathcal{Y} , and $\sigma^* \in \{\pm 1\}^n$, we say that $y \sim \mathsf{SI}(\sigma^*, \mathcal{S}_+, \mathcal{S}_-)$, where $y \in \mathcal{Y}^n$ such that for any $i \in [n]$, $y_i \sim \mathcal{S}_+, \text{ if } \sigma_i^* = +1 \text{ and } y_i \sim \mathcal{S}_- \text{ if } \sigma_i^* = -1.$

The entries $\{y_i\}_{i \in [n]}$ are independent conditional on σ^* . We assume that the likelihoods $S_+(y_i)$ and $\mathcal{S}_{-}(y_i)$ are computable in O(1) time.

Definition 2.5 (Gaussian Features (GF)). The model $y \sim \mathsf{GF}(\sigma^*, v_+, v_-, \sigma^2)$ for $v_+, v_- \in \mathbb{R}^d$ and $\sigma^2 > 0$ is a special case of SI with $\mathcal{Y} = \mathbb{R}^d$ and $\mathcal{S}_+ \equiv \mathcal{N}(v_+, \sigma^2 I_d)$ and $\mathcal{S}_- \equiv \mathcal{N}(v_-, \sigma^2 I_d)$.

Definition 2.6 (Binary Erasure Channel (BEC)). For any $\sigma^* \in \{\pm 1\}^n$ and $\epsilon \in (0, 1]$, we say $y \sim \mathsf{BEC}(\sigma^*, \epsilon)$ where each entry of σ^* is erased to 0, independently with probability ϵ , to form $y \in \{-1, 0, +1\}^n$.

Definition 2.7 (Binary Symmetric Channel (BSC)). For any $\sigma^* \in \{\pm 1\}^n$ and $\alpha \in (0, 1/2]$, we say $y \sim \mathsf{BSC}(\sigma^*, \alpha)$ where each entry of σ^* is flipped independently with probability α , to form $y \in \{\pm 1\}^n$.



Figure 1: Visualization of BSC and BEC side information. The red-colored, blue-colored, and uncolored vertices have side information labels of +1, -1, and 0 respectively.

2.2 EXACT RECOVERY.

Our goal is to exactly recover the community labels σ^* given the observation matrix A and the side information y when available, as formalized below.

Definition 2.8 (Exact Recovery). We say that an estimator $\hat{\sigma}$ succeeds if

(i) $\hat{\sigma} \in \{\pm \sigma^*\}$ when $\mathcal{P}_+ \equiv \mathcal{P}_-, \rho = 1/2$ and there is no side information (the symmetric

(ii) $\hat{\sigma} = \sigma^*$, when $\mathcal{P}_+ \not\equiv \mathcal{P}_-$ or $\rho \neq 1/2$ or side information is present (non-symmetric case).

Otherwise, we say $\hat{\sigma}$ fails. We say that $\hat{\sigma}$ achieves exact recovery if $\lim_{n \to \infty} \mathbb{P}(\hat{\sigma} \text{ succeeds}) = 1$.

Note that in the symmetric setting described in Definition 2.8, it is not possible to recover σ^* with high probability, and we can only hope to recover the partition. In all other cases, we wish to recover the labels and not just the partition. All our positive results will demonstrate recovery in this strong sense, while our negative results will show that even recovering the partition is impossible below the threshold. The optimal predictor for any model and side information is the one which has maximum posterior probability given the observation matrix A and the side information y, when it is present.

Definition 2.9 (MAP Estimator). Consider the observation matrix A and the side information y (either BSC or BEC). We define the Maximum A Posteriori (MAP) estimator as

$$\hat{\sigma}_{\text{MAP}} = \underset{\sigma \in \{\pm 1\}^n}{\arg \max} \mathbb{P}(\sigma^* = \sigma \mid A, y)$$

When no side information is present, define $\hat{\sigma}_{MAP} = \arg \max_{\sigma \in \{\pm 1\}^n} \mathbb{P}(\sigma^* = \sigma \mid A)$.

MAIN RESULTS

Before stating our algorithmic result, we describe the information-theoretic limit due to Dreveton et al. (2024) in our notation.

3.1 INFORMATION THEORETIC THRESHOLD FROM DREVETON ET AL. (2024)

272 Define the *Chernoff coefficient* across the pair of communities as

273 274 275

276

277

283 284

285

287

288 289

290

 $\mathbf{CH}_t(+,-) = (1-t) \left[\rho \mathbf{D}_t(\mathcal{P}_+ \| \mathcal{Q}) + (1-\rho) \mathbf{D}_t(\mathcal{Q} \| \mathcal{P}_-) + \frac{1}{n} \mathbf{D}_t(\mathcal{S}_+ \| \mathcal{S}_-) \right]$

Here $D_t(\mathcal{A}||\mathcal{B})$ is the st the Rényi divergence of order t between any two laws $(\mathcal{A}, \mathcal{B})$, such that they are both continuous or discrete, is given by

$$D_t(\mathcal{A}||\mathcal{B}) := \frac{1}{(t-1)} \log \mathbb{E}_{x \sim \mathcal{B}} \left[\left(\frac{\mathcal{A}(x)}{\mathcal{B}(x)} \right)^t \right].$$
(3)

Define the limit $L: (0,1) \to \mathbb{R}_{\geq 0} \cup \{+\infty\}$ by

$$L(t) = \lim_{n \to \infty} \frac{n}{\log n} CH_t(+, -).$$
(4)

We restrict ourselves to the laws such that L(t) is well defined. Then the information theoretic limit is characterized by

$$I^* = \sup_{t \in (0,1)} L(t),$$
 (5)

where by convention, we consider the supremum to be $+\infty$ if L(t) is unbounded in $t \in (0, 1)$ or there exists $t \in (0, 1)$ such that $L(t) = +\infty$. Then following proposition is a minor variant of the two community case of (Dreveton et al., 2024, Theorem 1), characterizing the fundamental limit for exact recovery for the GBM.

Proposition 3.1 (Dreveton et al. (2024)). Let $\rho \in (0,1)$ and $(\mathcal{P}_+, \mathcal{P}_-, \mathcal{Q})$ be probability laws. Let $(A, \sigma^*) \sim \text{GBM}_n(\rho, \mathcal{P}_+, \mathcal{P}_-, \mathcal{Q})$ and we observe A. Optionally, for side information laws $(\mathcal{S}_+, \mathcal{S}_-)$, we observe $y \sim \text{SI}(\sigma^*, \mathcal{S}_+, \mathcal{S}_-)$. We assume the laws are such that I^* in (5) is welldefined. Then

1. If $I^* > 1$, then the Maximum A Posteriori estimator $\hat{\sigma}_{MAP}$ achieves exact recovery.

300 301 302

303

2. Additionally, if $I^* < 1$ and the function L(t) is strictly concave then no estimator $\hat{\sigma}$ achieves exact recovery.

This variant follows from (Dreveton et al., 2024, Theorem 1) by observing that the success of MAP proof does not rely on the strict concavity of L(t) and the proof goes through even if $I^* = +\infty$. The impossibility result is completely equivalent to theirs. We discuss the necessary proof changes in Appendix A.3. We note that the case of no side information can be realized by considering side information laws (S_+, S_-) that always deterministically output 0; in this case $D_t(S_+||S_-) = 0$.

309 310

3.2 Algorithmic Achievability for ROS and SBM

The main results of the paper is to design an optimal spectral algorithm for ROS and SBM.

Theorem 1. Fix $\rho \in (0, 1)$ and $a, b \in \mathbb{R}$ such that $\max\{|a|, |b|\} > 0$. Let $(A, \sigma^*) \sim \mathsf{ROS}_n(\rho, a, b)$. Optionally, consider a side information channel $y \sim \mathsf{SI}(\sigma^*, \mathcal{S}_+, \mathcal{S}_-)$ such that I^* in (5) is welldefined. Then there is a spectral algorithm (Algorithm 2) that returns the estimator $\hat{\sigma}_{spec}$ which achieves exact recovery, whenever $I^* > 1$.

Theorem 2. Let $\rho \in (0,1)$ and $a_1, a_2, b > 0$. Let $(A, \sigma^*) \sim \text{SBM}_n(\rho, a_1, a_2, b)$. Optionally, consider a side information channel $y \sim \text{SI}(\sigma^*, S_+, S_-)$ such that I^* in (5) is well-defined. Then there is a spectral algorithm (Algorithms 5 and 6) that returns the estimator $\hat{\sigma}_{\text{spec}}$ which achieves exact recovery, whenever $I^* > 1$.

In Appendix C, we will verify that for both SBM and ROS, the function L(t) (and thus also I^*) is well-defined and L(t) is strictly concave for each (i) no side information (ii) GF, BEC, and BSC channels. This along with Theorems 1, 2 and Proposition 3.1 establishes that the spectral algorithm succeeds up to the information-theoretic limit in these settings.

324 4 GENIE-AIDED ESTIMATORS

 $\hat{\sigma}$

We begin our analysis by defining the framework of genie-aided estimation (see e.g., Abbe (2017)). Our contribution is a systematic method of connecting spectral algorithms to genie estimators. In the genie-aided setting, we suppose that all labels but the i^{th} are known, and the goal is to determine the i^{th} label. More formally, let σ_{-i}^* denote the true labels, apart from σ_i^* . The optimal estimator for the i^{th} label is given by

331

$$\mathcal{G}_{\text{Gen},i} = \begin{cases} +1, & \text{if } \mathbb{P}(\sigma_i^* = +1 \mid A, y, \sigma_{-i}^*) \ge \mathbb{P}(\sigma_i^* = -1 \mid A, y, \sigma_{-i}^*).\\ -1, & \text{otherwise.} \end{cases}$$
(6)

(where the conditioning on y is omitted when there is no side information present). Moreover, we say that $\hat{\sigma}_{\text{Gen},i}$ fails on an instance if the posterior probability of the incorrect label is *strictly* greater than the posterior probability of the correct label. The following lemma rigorously establishes the intuitive claim that the failure of some genie-aided estimator implies the global MAP also fails.

338 Lemma 4.1. Let $\rho \in (0,1)$ and $\mathcal{P}_+, \mathcal{P}_-, \mathcal{Q}$ be any distributions. Let $(A, \sigma^*) \sim GBM_n(\rho, \mathcal{P}_+, \mathcal{P}_-, \mathcal{Q})$. Optionally, let $y \sim SI(\sigma^*, S_+, S_-)$ in \mathcal{Y}^n where $y \in \mathcal{Y}^n$ is the side information for any laws (S_+, S_-) over \mathcal{Y} . Define the genie-aided estimators $\{\hat{\sigma}_{Gen,i} : i \in [n]\}$ and $\hat{\sigma}_{MAP}$ for the respective model. Then

$$\mathbb{P}(\exists i \in [n] : \hat{\sigma}_{\text{Gen},i} \text{ fails}) \leq \mathbb{P}(\hat{\sigma}_{\text{MAP}} \text{ fails})$$

Genie scores. We form a vector $z^* \in (\mathbb{R} \cup \{\pm \infty\})^n$, called the *genie score* vector, where z_i^* records the log of the ratio of posterior probabilities of the label σ_i^* given σ_{-i}^* by a genie. That is,

342 343 344

 $z_{i}^{*} = \log \left(\frac{\mathbb{P}(\sigma_{i}^{*} = +1 \mid A, y, \sigma_{-i}^{*})}{\mathbb{P}(\sigma_{i}^{*} = -1 \mid A, y, \sigma_{-i}^{*})} \right) \quad \text{or} \quad z_{i}^{*} = \log \left(\frac{\mathbb{P}(\sigma_{i}^{*} = +1 \mid A, \sigma_{-i}^{*})}{\mathbb{P}(\sigma_{i}^{*} = -1 \mid A, \sigma_{-i}^{*})} \right), \tag{7}$

in the cases of side information and no side information respectively. Then the optimal geniebased estimator corresponds to $\hat{\sigma}_{\text{Gen},i} = \text{sgn}(z_i^*)$. Throughout the paper, we treat log as function from $[0,\infty]$ to the extended real line and implicitly use standard conventions of extended real line algebra. See Appendix B.1. The following lemma gives the form of the genie scores without or with node-attributed side information. Here we define $C_{+}^{-i} := C_{+} \setminus \{i\}$ and $C_{-}^{-i} := C_{-} \setminus \{i\}$.

Lemma 4.2. Consider any $\rho \in (0,1)$ and laws $\mathcal{P}_+, \mathcal{P}_-$ and \mathcal{Q} . Let $(A, \sigma^*) \sim \mathsf{GBM}_n(\rho, \mathcal{P}_+, \mathcal{P}_-, \mathcal{Q})$. Optionally, let $y \sim \mathsf{SI}(\sigma^*, \mathcal{S}_+, \mathcal{S}_-)$ where $y \in \mathcal{Y}^n$ is the side information for any laws $(\mathcal{S}_+, \mathcal{S}_-)$ over \mathcal{Y} . Then for any $i \in [n]$, the genie score for the *i*th label is given by

• No side information:

$$z_i^* = \sum_{j \in C_+^{-i}} \log\left(\frac{\mathcal{P}_+(A_{ij})}{\mathcal{Q}(A_{ij})}\right) + \sum_{j \in C_-^{-i}} \log\left(\frac{\mathcal{Q}(A_{ij})}{\mathcal{P}_-(A_{ij})}\right) + \log\left(\frac{\rho}{1-\rho}\right).$$

 $z_i^* = \sum_{j \in C_+^{-i}} \log\left(\frac{\mathcal{P}_+(A_{ij})}{\mathcal{Q}(A_{ij})}\right) + \sum_{j \in C^{-i}} \log\left(\frac{\mathcal{Q}(A_{ij})}{\mathcal{P}_-(A_{ij})}\right) + \log\left(\frac{\rho}{1-\rho}\right) + \log\left(\frac{\mathcal{S}_+(y_i)}{\mathcal{S}_-(y_i)}\right).$

• Side information:

354 355

356 357

368 Due to the independence of A and y, conditioned on σ^* , the genie score z_i^* under side information is simply the score without side information with the addition of the log-likelihood ratio of the node 369 attribute y_i . Intuitively, we devise a principled method for determining the weights of the eigenvector 370 combination and threshold value in our spectral algorithm such that ith entry of the vector formed 371 approximates the the genie score z_i^* under no side information. Observe that the statistic z_i^* only 372 depends on the i^{th} row of A. Remarkably, the spectral algorithm will approximate this statistic from 373 eigenvectors of A, despite not knowing C_{+}^{i} and C_{-}^{i} . In light of Lemma 4.2, it is also immediate 374 to see the motivation behind the use of side information prescribed in Section 1.1. In particular, the 375 genie scores undergo exactly the same transformation when the side information becomes available 376 according to Lemma 4.2. 377

Lemma 4.3. For any $i \in [n]$, the log-likelihood ratio of side information label is given by

• <u>GF (Gaussian features)</u>: When $y \sim GF(\sigma^*, v_+, v_-, \sigma^2)$ for $v_+, v_- \in \mathbb{R}^d$ and $\sigma^2 > 0$,

$$\log\left(\frac{S_{+}(y_{i})}{S_{-}(y_{i})}\right) = \frac{\|y_{i} - v_{-}\|_{2}^{2} - \|y_{i} - v_{+}\|_{2}^{2}}{2\sigma^{2}}$$

• BEC channel: When $y \sim \text{BEC}(\sigma^*, \epsilon)$ for $\epsilon \in (0, 1]$,

$$\log\left(\frac{\mathcal{S}_{+}(y_{i})}{\mathcal{S}_{-}(y_{i})}\right) = \begin{cases} +\infty, & \text{if } \sigma_{i}^{*} = +1; \\ -\infty, & \text{if } \sigma_{i}^{*} = -1; \\ 0, & \text{otherwise.} \end{cases}$$

• BSC channel: When $y \sim BSC(\sigma^*, \alpha)$ for $\alpha \in (0, 0.5]$,

$$\log\left(\frac{\mathcal{S}_{+}(y_{i})}{\mathcal{S}_{-}(y_{i})}\right) = \log\left(\frac{1-\alpha}{\alpha}\right)y_{i}$$

Genie Success with Margin above the IT Limit: We start by a crucial observation about the genie scores, which will be important in analyzing the spectral algorithm. We show that above the information-theoretic limit, the genie score have the correct sign corresponding to the true label, but with a *sufficient margin*. We show the for any model GBM and optionally a side information channel, whenever $I^* > 1$, there exists a constant $\delta > 0$ such that with high probability

$$\min_{i \in [n]} \sigma_i^* z_i^* > \delta \log n, \tag{8}$$

formalized in Lemma D.1.

5 GENIE TO SPECTRAL ALGORITHM (FOR ROS AND SBM)

In light of (8), if an algorithm can approximate the genie score up to an additive error of $o(\log n)$, then sign thresholding will correctly recover all the labels. We first establish that for both ROS and SBM the genie score vector takes a special form

403 404 405

411

success).

378

379380381382

384

386 387

388 389

390 391

392

393 394

395

396 397

398 399

400 401

402

$$z^* \approx Aw + \gamma \mathbf{1}_n,$$

where the approximation is in the ℓ_{∞} norm, for a certain $\gamma \in \mathbb{R}$ and $w \in \mathbb{R}^n$ with entries $(w_+, w_-) \in \mathbb{R}^2$ in the locations of C_+ and C_- respectively. See Lemmas F.1 (ROS) and G.2 (SBM) for precise statements. Note that (w_+, w_-, γ) are just scalars that can be calculated from the model parameters and do not depend on σ^* . The main power of the genie lies in forming the vector w, which requires knowing the locations C_+ and C_- . Then the question that remains is how one may come up with a proxy for w such that the genie score is well-approximated.

412 5.1 WARM-UP: DEGREE PROFILING ALGORITHM FOR BEC AND BSC CHANNELS 413

Under BEC and BSC channel, a natural question is then what if we simply trust the side information labels $y \in \{-1, 0, +1\}^n$ and $y \in \{-1, +1\}^n$ respectively. In particular, we use the proxy for wwhere we just trust the side information on the face value and use the locations of $S_+ := \{i \in [n] : y_i = +1\}$ and $S_- := \{i \in [n] : y_i = -1\}$ instead of C_+ and C_- .

It is known that the asymptotic information-theoretic threshold does not shift for the BEC and BSC channels unless $\epsilon = O(n^{-\beta})$ and $\alpha = O(n^{-\beta})$ for some $\beta > 0$ (Dreveton et al., 2024; Saad & Nosratinia, 2018). Therefore, the side information y already satisfies almost exact recovery criterion, recovering (1 - o(1)) labels correctly. Hence, we have $|C_+\Delta S_+|, |C_-\Delta S_-| = o(n)$ with high probability. We then show that for this proxy choice of scores z^{dp} , we indeed have $||z^{dp} - z^*||_{\infty} = o(\log n)$, and thus, the degree-profiling algorithm succeeds down to the shifted threshold. The formal theorems and algorithms can be found in Appendix F.3 and G.3.

424 **Remark 5.1.** We emphasize that the degree profiling algorithm has an important caveat that it 425 would fail to recover labels exactly from a tuple (A, y), when side information strength is "weak" 426 or completely absent, even though the recovery was possible just from A. To overcome this, one has 427 to rely on the signal from A to get preliminary almost exactly correct labels, rather than just trusting 428 the side information y. This exactly corresponds to the two-stage strategies described in Section 1.2. 429 Our spectral algorithm in just one stage recover all the labels correctly from (A, y) (including no side information), whenever it is possible to do so, and that too for any side information channel 430 (under the technical assumption that I^* in (5) is well-defined which is also needed for the MAP 431

432 5.2 SPECTRAL ALGORITHM

Spectral Algorithm. The spectral algorithm affords a significantly more versatility than the 435 degree-profiling algorithm, as it emulates the genie score z^* without any clean-up step with more 436 general or even no side information. The design of our spectral algorithm is informed by the en-437 trywise eigenvector analysis result of Abbe et al. (2020), which allows us to say that the leading K438 eigenvectors of A satisfy

$$u_i \approx \frac{Au_i^*}{\lambda_i^*},$$

where (λ_i^*, u_i^*) is the corresponding eigenvector of the expectation matrix $\mathbb{E}[A \mid \sigma^*]$, and the approximation is in the ℓ_{∞} norm. Here K = 1 for ROS and K = 2 for SBM. For both models, the matrix $\mathbb{E}[A \mid \sigma^*]$ has a block structure, so do its top eigenvectors as well. Thus, we can find an appropriate linear combination coefficients $(c_i)_{i \in [K]}$ such that

$$w \approx \sum_{i=1}^{K} \frac{c_i}{\lambda_i^*} u_i^*.$$

448 It is important to note that computing $(c_i)_{i \in [K]}$ does not require knowing the ground-truth σ^* .

Algorithm 1 An informal sketch of the spectral algorithm

Input: An observation matrix $A \in \mathbb{R}^{n \times n}$ and the model parameters. Optionally, the side information $y \in \mathcal{Y}^n$.

Output: An estimate of community assignments $\hat{\sigma}_{\text{spec}}$.

1: Compute coefficients $(c_i)_{i \in [K]}$ from the model parameters such that $w \approx \sum_{i=1}^{K} \frac{c_i}{\lambda^*} u_i^*$.

2: Compute the top K eigenpairs of A denoted by (u_i, λ_i) .

3: Form the spectral score vector:

• No side information:

$$z^{\text{spec}} = \gamma \mathbf{1}_n + \sum_{i=1}^K c_i u_i.$$

• Side information: If side information is present, further update

$$z^{\text{spec}} = z^{\text{spec}} + \log\left(\frac{\mathcal{S}_+}{\mathcal{S}_-}(y)\right).$$

4: $\hat{\sigma}_{\text{spec}} = \text{sgn}(z^{\text{spec}}).$

Then the spectral score vector when the side information is absent given by

$$z^{\text{spec}} = \gamma \mathbf{1}_n + \sum_{i=1}^K c_i u_i \approx \gamma \mathbf{1}_n + \sum_{i=1}^K c_i \frac{Au_i^*}{\lambda_i^*} \approx \gamma \mathbf{1}_n + Aw \approx z^*.$$

This achieves an ℓ_{∞} approximation of the genie score under no side information using eigenvectors of A. When side information is present, due to Lemma 4.2, it suffices to add the log-likelihood ratio vector from (1). The approximations are tight enough such that we achieve

$$\left\| z^{\text{spec}} - z^* \right\|_{\infty} = o(\log n).$$

478 Recalling that the genie scores succeed with margin of $\Omega(\log n)$ by (8) above the IT threshold 479 immediately gives us optimality of spectral algorithm.

Remark 5.2. We remark that, to show the impossibility of recovery, due to the optimality of MAP,
Dreveton et al. (2024) shows that the MAP estimator (Definition 2.9) fails below the threshold.
However, their proof indeed shows even stronger statement that below the threshold, even the genieaided estimators fail. Both their impossibility and our achievability of spectral algorithm results
are driven by the genie-aided estimators, so it comes as no surprise that there is a certain threshold collapse: the genie-aided estimators, spectral estimator, and MAP estimator all achieve the same recovery threshold. The entire discussion can be summarized in Figure 2.



Figure 2: Summary of our unified proof framework. Our spectral algorithm is designed such that (*a*) holds. Note that: (*b*) follows from the optimality of $\hat{\sigma}_{MAP}$ and (*c*) follows from Lemma 4.1. Finally, we get the threshold collapse phenomenon because below the IT threshold, the event in the fourth block happens with probability o(1), and above the IT threshold, the event in the first block happens with high probability (8).

6 DISCUSSION AND FUTURE WORK

490

491

492

493

494 495

496

497

498

499

500

501

504

505

506

507

510

511

522

523

524 525

526 527

528

529

In this paper, we provide a systematic treatment of designing optimal spectral algorithms for twocommunity matrix inference problems under side information, focused on the Bernoulli and Gaussian cases. From a technical standpoint, our work makes a rigorous connection between spectral algorithms and genie-aided estimators, characterizing their effectiveness in achieving sharp thresholds for various exact recovery problems in a recent line of work. We refer the reader to Appendix A.1 for a detailed discussion on this. Understanding the capabilities of such vanilla spectral algorithms, without any clean-up stage, is of fundamental interest; we hope this perspective will guide the design and analysis of spectral algorithms for exact community recovery problems moving forward. Some directions for future work include:

- Exact recovery in Gaussian Mixture Block Model: In a recent work, Li & Schramm (2023) proposed an alternative model to better capture real-world networks and sketched out the general landscape for recovery by studying almost exact recovery. What about exact recovery? Interestingly, Li & Schramm (2023) proposed exactly the same vanilla spectral algorithm and showed it achieves almost exact recovery. Does it also succeed for exact recovery?
- 512 • More general degree-profiling: A natural analogue of degree-profiling algorithm for more 513 general side information beyond BEC and BSC is to sign threshold the log-likelihood ratio 514 vector $\log\left(\frac{S_+}{S}(y)\right)$ to compute preliminary labeling. We conjecture that, whenever such 515 side information is sufficient to shift the exact recovery threshold, the preliminary assign-516 ment will already correctly compute (1 - o(1))-faction of labels. Emulating the genie then 517 using this labeling as a proxy should succeed in exact recovery. We again emphasize that 518 the degree-profiling algorithm has limitations (Remark 5.1), however, it can be seen as a 519 good alternative of the spectral algorithms in the limited scope when side information is guaranteed to be enormous. 521
 - More general settings: To design spectral algorithms for more than two communities and more general observation distributions $(\mathcal{P}_+, \mathcal{P}_-, \mathcal{Q})$ from a class of exponential families (also see details in Appendix A.1).

References

- Emmanuel Abbe. Community detection and stochastic block models: recent developments. *The Journal of Machine Learning Research*, 18(1):6446–6531, 2017.
- Emmanuel Abbe, Afonso S Bandeira, Annina Bracher, and Amit Singer. Decoding binary node
 labels from censored edge measurements: Phase transition and efficient recovery. *IEEE Transac- tions on Network Science and Engineering*, 1(1):10–22, 2014.
- Emmanuel Abbe, Afonso S. Bandeira, and G. Hall. Exact recovery in the stochastic block model. *IEEE Transactions on Information Theory*, 62:471–487, 2016.
- Emmanuel Abbe, Jianqing Fan, Kaizheng Wang, and Yiqiao Zhong. Entrywise eigenvector analysis
 of random matrices with low expected rank. *Annals of Statistics*, 48 3:1452–1474, 2020.
- 539 Emmanuel Abbe, Francois Baccelli, and Abishek Sankararaman. Community detection on euclidean random graphs. *Information and Inference: A Journal of the IMA*, 10(1):109–160, 2021.

- 540 Noga Alon, Michael Krivelevich, and Benny Sudakov. Finding a large hidden clique in a random 541 graph. Random Structures & Algorithms, 13(3-4):457–466, 1998. 542 Eric Bair. Semi-supervised clustering methods. Wiley Interdisciplinary Reviews: Computational 543 Statistics, 5(5):349–361, 2013. 544 Ramnath Balasubramanyan and William W Cohen. Block-LDA: Jointly modeling entity-annotated 546 text and entity-entity links. In Proceedings of the 2011 SIAM International Conference on Data 547 Mining, pp. 450-461. SIAM, 2011. 548 Afonso S Bandeira, Nicolas Boumal, and Amit Singer. Tightness of the maximum likelihood 549 semidefinite relaxation for angular synchronization. *Mathematical Programming*, 163:145–167, 550 2017. 551 552 Sugato Basu, Mikhail Bilenko, and Raymond J Mooney. A probabilistic framework for semi-553 supervised clustering. In Proceedings of the tenth ACM SIGKDD international conference on 554 Knowledge discovery and data mining, pp. 59–68, 2004. 555 Norbert Binkiewicz, Joshua T Vogelstein, and Karl Rohe. Covariate-assisted spectral clustering. 556 Biometrika, 104(2):361-377, 2017. 558 Cécile Bothorel, Juan David Cruz, Matteo Magnani, and Barbora Micenkova. Clustering attributed 559 graphs: models, measures and methods. Network Science, 3(3):408-444, 2015. Guillaume Braun, Hemant Tyagi, and Christophe Biernacki. An iterative clustering algorithm for 561 the contextual stochastic block model with optimality guarantees. In International Conference on 562 Machine Learning, pp. 2257–2291. PMLR, 2022. 563 564 Guy Bresler, Chenghao Guo, and Yury Polyanskiy. Thresholds for reconstruction of random hyper-565 graphs from graph projections. In The Thirty Seventh Annual Conference on Learning Theory, pp. 632-647. PMLR, 2024. 566 567 T Tony Cai, Tengyuan Liang, and Alexander Rakhlin. Inference via message passing on partially 568 labeled stochastic block models. arXiv preprint arXiv:1603.06923, 2016. 569 570 T Tony Cai, Tengyuan Liang, and Alexander Rakhlin. Computational and statistical boundaries for submatrix localization in a large noisy matrix. The Annals of Statistics, 45(4):1403–1430, 2017. 571 572 Jonathan Chang and David M Blei. Hierarchical relational models for document networks. The 573 Annals of Applied Statistics, pp. 124–150, 2010. 574 575 Olivier Chapelle, Jason Weston, and Bernhard Schölkopf. Cluster kernels for semi-supervised learning. Advances in neural information processing systems, 15, 2002. 576 577 Hong Cheng, Yang Zhou, and Jeffrey Xu Yu. Clustering large attributed graphs: A balance be-578 tween structural and attribute similarities. ACM Transactions on Knowledge Discovery from Data 579 (*TKDD*), 5(2):1–33, 2011. 580 Amin Coja-Oghlan. A spectral heuristic for bisecting random graphs. Random Structures & Algo-581 rithms, 29(3):351-398, 2006. 582 583 Aurelien Decelle, Florent Krzakala, Cristopher Moore, and Lenka Zdeborová. Asymptotic analysis 584 of the stochastic block model for modular networks and its algorithmic applications. *Physical* 585 *Review E*, 84(6):066106, 2011. 586 Shaofeng Deng, Shuyang Ling, and Thomas Strohmer. Strong consistency, graph Laplacians, and the stochastic block model. J. Mach. Learn. Res., 22(117):1-44, 2021. 588 589 Yash Deshpande, Subhabrata Sen, Andrea Montanari, and Elchanan Mossel. Contextual stochastic block models. Advances in Neural Information Processing Systems, 31, 2018. Souvik Dhara, Julia Gaudio, Elchanan Mossel, and Colin Sandon. Spectral recovery of binary 592
- Souvik Dhara, Julia Gaudio, Elchanan Mossel, and Colin Sandon. Spectral recovery of binary censored block models. In *Proceedings of the 2022 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 3389–3416. SIAM, 2022a.

- 594 Souvik Dhara, Julia Gaudio, Elchanan Mossel, and Colin Sandon. Spectral algorithms optimally recover (censored) planted dense subgraphs. arXiv preprint arXiv:2203.11847, 2022b. 596 Souvik Dhara, Julia Gaudio, Elchanan Mossel, and Colin Sandon. The power of two matrices in 597 spectral algorithms for community recovery. IEEE Transactions on Information Theory, 2023. 598 Maximilien Dreveton, Felipe Fernandes, and Daniel Figueiredo. Exact recovery and bregman hard 600 clustering of node-attributed stochastic block model. Advances in Neural Information Processing 601 Systems, 36, 2024. 602 Julia Gaudio and Nirmit Joshi. Community detection in the hypergraph SBM: Exact recovery given 603 the similarity matrix. In The Thirty Sixth Annual Conference on Learning Theory, pp. 469–510. 604 PMLR, 2023. 605 Julia Gaudio and Heming Liu. Spectral recovery in the labeled sbm. arXiv preprint 607 arXiv:2408.13075, 2024. 608 Julia Gaudio, Miklos Z Racz, and Anirudh Sridhar. Exact community recovery in correlated stochas-609 tic block models. In Conference on Learning Theory, pp. 2183–2241. PMLR, 2022. 610 611 Julia Gaudio, Charlie Guan, Xiaochun Niu, and Ermin Wei. Exact label recovery in euclidean 612 random graphs. arXiv preprint arXiv:2407.11163, 2024a. 613 614 Julia Gaudio, Xiaochun Niu, and Ermin Wei. Exact community recovery in the geometric sbm. 615 In Proceedings of the 2024 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA), pp. 2158-2184. SIAM, 2024b. 616 617 Jaume Gibert, Ernest Valveny, and Horst Bunke. Graph embedding in vector spaces by node attribute 618 statistics. *Pattern Recognition*, 45(9):3072–3083, 2012. 619 620 Manuel Gil, Fady Alajaji, and Tamas Linder. Rényi divergence measures for commonly used uni-621 variate continuous distributions. Information Sciences, 249:124–131, 2013. 622 Stephan Günnemann, Ines Färber, Sebastian Raubach, and Thomas Seidl. Spectral subspace cluster-623 ing for graphs with feature vectors. In 2013 IEEE 13th International Conference on Data Mining, 624 pp. 231–240. IEEE, 2013. 625 626 Bruce Hajek, Yihong Wu, and Jiaming Xu. Exact recovery threshold in the binary censored block 627 model. In 2015 IEEE Information Theory Workshop-Fall (ITW), pp. 99–103. IEEE, 2015. 628 Bruce Hajek, Yihong Wu, and Jiaming Xu. Achieving exact cluster recovery threshold via semidef-629 inite programming. IEEE Transactions on Information Theory, 62(5):2788-2797, 2016. 630 631 Bruce Hajek, Yihong Wu, and Jiaming Xu. Submatrix localization via message passing. Journal of 632 *Machine Learning Research*, 18(186):1–52, 2018. 633 Adel Javanmard, Andrea Montanari, and Federico Ricci-Tersenghi. Phase transitions in semidefinite 634 relaxations. Proceedings of the National Academy of Sciences, 113(16):E2218-E2223, 2016. 635 636 Yicong Jiang and Tracy Ke. Semi-supervised community detection via structural similarity metrics. 637 In The Eleventh International Conference on Learning Representations, 2023. 638 Arun Kadavankandy, Konstantin Avrachenkov, Laura Cottatellucci, and Rajesh Sundaresan. The 639 power of side-information in subgraph detection. IEEE Transactions on Signal Processing, 66 640 (7):1905–1919, 2017. 641 642 Varun Kanade, Elchanan Mossel, and Tselil Schramm. Global and local information in clustering 643 labeled block models. *IEEE Transactions on Information Theory*, 62(10):5906–5917, 2016. 644 Anastasia Krithara, Massih R Amini, Jean-Michel Renders, and Cyril Goutte. Semi-supervised 645 document classification with a mislabeling error model. In Advances in Information Retrieval: 646
- document classification with a mislabeling error model. In Advances in Information Retrieval:
 30th European Conference on IR Research, ECIR 2008, Glasgow, UK, March 30-April 3, 2008.
 Proceedings 30, pp. 370–381. Springer, 2008.

- Shuangping Li and Tselil Schramm. Spectral clustering in the gaussian mixture block model. *arXiv* preprint arXiv:2305.00979, 2023.
- Frank McSherry. Spectral partitioning of random graphs. In *Proceedings 42nd IEEE Symposium on Foundations of Computer Science*, pp. 529–537. IEEE, 2001.
- Michael Mitzenmacher and Eli Upfal. Probability and computing: Randomization and probabilistic techniques in algorithms and data analysis. Cambridge university press, 2017.
- Javad Zahedi Moghaddam, Mohammad Esmaeili, and Aria Nosratinia. Exact recovery threshold in
 dynamic binary censored block model. In 2022 IEEE International Symposium on Information
 Theory (ISIT), pp. 1088–1093. IEEE, 2022.
- Elchanan Mossel and Jiaming Xu. Local algorithms for block models with side information. In Proceedings of the 2016 ACM Conference on Innovations in Theoretical Computer Science, pp. 71–80, 2016.
- 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*, pp. 69–75, 2015.
- Mark EJ Newman and Aaron Clauset. Structure and inference in annotated networks. *Nature communications*, 7(1):11863, 2016.
- Li Ni, Junnan Ge, Yiwen Zhang, Wenjian Luo, and Victor S. Sheng. Semi-supervised local commu nity detection. *IEEE Transactions on Knowledge and Data Engineering*, 36(2):823–839, 2024.
 doi: 10.1109/TKDE.2023.3290095.
- Miklos Racz and Anirudh Sridhar. Correlated stochastic block models: Exact graph matching with applications to recovering communities. *Advances in Neural Information Processing Systems*, 34: 22259–22273, 2021.
- Hussein Saad and Aria Nosratinia. Community detection with side information: Exact recovery under the stochastic block model. *IEEE Journal of Selected Topics in Signal Processing*, 12(5): 944–958, 2018.
- Hussein Saad and Aria Nosratinia. Recovering a single community with side information. *IEEE Transactions on Information Theory*, 66(12):7939–7966, 2020.

681

- Jesper E Van Engelen and Holger H Hoos. A survey on semi-supervised learning. *Machine Learning*, 109(2):373–440, 2020.
- Van Vu. A simple svd algorithm for finding hidden partitions. *Combinatorics, Probability and Computing*, 27(1):124–140, 2018.
- Xixi Wu, Yun Xiong, Yao Zhang, Yizhu Jiao, Caihua Shan, Yiheng Sun, Yangyong Zhu, and Philip S
 Yu. Clare: A semi-supervised community detection algorithm. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*, pp. 2059–2069, 2022.
- ⁶⁹⁰
 ⁶⁹¹ Zhiqiang Xu, Yiping Ke, Yi Wang, Hong Cheng, and James Cheng. A model-based approach to attributed graph clustering. In *Proceedings of the 2012 ACM SIGMOD international conference on management of data*, pp. 505–516, 2012.
- Jaewon Yang, Julian McAuley, and Jure Leskovec. Community detection in networks with node attributes. In 2013 IEEE 13th international conference on data mining, pp. 1151–1156. IEEE, 2013.
- Tianbao Yang, Rong Jin, Yun Chi, and Shenghuo Zhu. Combining link and content for community detection: a discriminative approach. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 927–936, 2009.
- 701 Se-Young Yun and Alexandre Proutiere. Accurate community detection in the stochastic block model via spectral algorithms. *arXiv preprint arXiv:1412.7335*, 2014.

Se-Young Yun and Alexandre Proutiere. Optimal cluster recovery in the labeled stochastic block model. Advances in Neural Information Processing Systems, 29, 2016.

Hugo Zanghi, Stevenn Volant, and Christophe Ambroise. Clustering based on random graph model embedding vertex features. Pattern Recognition Letters, 31(9):830-836, 2010.

Ning Zhang, Weina Wang, and Lele Wang. Attributed graph alignment. In 2021 IEEE International Symposium on Information Theory (ISIT), pp. 1829–1834. IEEE, 2021.

Pan Zhang, Cristopher Moore, and Lenka Zdeborová. Phase transitions in semisupervised clustering of sparse networks. Physical Review E, 90(5):052802, 2014.

Yuan Zhang, Elizaveta Levina, and Ji Zhu. Community detection in networks with node features. Electronic Journal of Statistics, 10:3153-3178, 2016.

Yang Zhou, Hong Cheng, and Jeffrey Xu Yu. Graph clustering based on structural/attribute similarities. Proceedings of the VLDB Endowment, 2(1):718-729, 2009.

Appendices

723	Α	Defe	erred Discussions	15
724		A.1	Future Extensions and Prior Works on Spectral Algorithms	15
725		A.2	Additional Related Work	15
726		A.3	Proof Review of Dreveton et al. (2024) for Proposition 3.1	16
727		1110		10
728	В	Nota	ation and Probability Facts	17
730		B .1	Notation	17
731		B.2	Standard Probability Lemmas	17
732				
733	С	Veri	fication for Information Theoretic Limits	18
734	_			
735	D	Omi	itted Proofs from Section 4	20
736		D.1	Genie Success with Margin above the IT Threshold	22
738	Е	Ent	www.e. Rehavior of Figenvectors	23
739	Ľ	E 1	Entrusise Analysis for Eigenvectors for DOS	24
740		E.1	Entry wise Analysis for Eigenvectors for KOS	24
741		E.2	Entrywise Analysis of Eigenvectors for SBM	26
742	F	Proc	ofs and Algorithms for ROS.	28
743		F.1	Genie scores' formula when no side information	28
744		F.2	Spectral Algorithm and Proof of Theorem 1	29
746		F.3	Degree-Profiling Algorithm for BEC and BSC Channels	32
747		110		
748	G	Proc	ofs and Algorithms for SBM.	36
749		G.1	Genie scores' formula when no side information	36
750		G.2	Spectral Algorithm and Proof of Theorem 2	38
751		G 3	Degree-Profiling Algorithm for BEC and BSC Channels	41
752		0.5		.1
753				

756 A DEFERRED DISCUSSIONS757

758 A.1 FUTURE EXTENSIONS AND PRIOR WORKS ON SPECTRAL ALGORITHMS

759 Future Extensions. We now describe, how we can generalize the spectral algorithm to more than 760 two communities and more general distribution families. The entire framework of genie-aided esti-761 mation naturally generalizes to the multi-community case. We note that the entrywise eigenvectors 762 behaviors of Abbe et al. (2020) for top eigenvectors continue to hold. Despite this, (Abbe et al., 763 2020, Appendix C.4) noted difficulties in designing spectral algorithms for more than two blocks 764 due to the multiplicity of eigenvalues and eigenvectors being only computed up to a rotation. But 765 we note that, except for these degenerate cases (some measure zero subset of parameters where ex-766 act recovery is possible), the algorithm should be able to emulate the genie-aided estimation and achieve optimality. For the degenerate cases noted in Abbe et al. (2020), it remains an interesting 767 open question if we can design a spectral estimator. 768

We expect that more general distribution families for pairwise observation matrix could be possible
to handle. In particular, we first summarize the key technical properties we relied on for designing
the spectral algorithms for ROS and SBM.

- 1. The genie-score vector takes a special form where it is an affine function of the entries of the matrices.
 - $z^* \approx Aw + \gamma \mathbf{1}_n.$
- 2. The optimal genie linear combination vector w is in the span of eigenvectors of $\mathbb{E}[A \mid \sigma^*]$.
- 3. Lastly, the eigenvector u of A, and its corresponding eigenvector u^* of $\mathbb{E}[A \mid \sigma^*]$, by Abbe et al. (2016) we have
 - $u \approx \frac{Au^*}{\lambda^*},$

which allows us to emulate the genie using eigenvectors of A.

Property 1 holds more generally for any exponential family of distributions. The linear combinations in Property 2 needs to be designed for specific distribution families, but we expect this to be possible in great generality due to the block structure of the expectation matrices. Finally, one needs to verify the entrywise behavior of eigenvector *A* holds. Abbe et al. has general conditions under which this behavior holds, however, verifying the technical conditions for the distribution families remains to be the main challenge.

Prior Works. We note that prior works on spectral algorithms mentioned in Section 1 are all essentially emulating the genie and the model satisfies the aforementioned three properties. For example, prior works have used weighted adjacency matrices (Dhara et al., 2022a), and even multiple adjacency matrices (Dhara et al., 2023; Gaudio & Liu, 2024) in order to bring out "helpful" eigenvectors that can be used to mimic the genie estimators.

Interestingly, for the exact recovery problem in the hypergraph SBM from the similarity matrix W794 (counting the hyperedges involving each pair of vertices), Gaudio & Joshi (2023) devised the same 795 spectral strategy of computing the signs of the second eigenvector in the symmetric case and ana-796 lyzed it by proving the ℓ_{∞} behavior of eigenvectors. However, this strategy is known to be strictly 797 suboptimal (Bresler et al., 2024) for this setting. This is because the genie scores are not affine 798 functions of the entries of the matrices W due to internal dependencies. However, it is interesting 799 to note that the algorithm achieves the threshold of the min-bisection estimator (also not efficiently 800 computable), which is the first order approximation of MAP. In other words, the eigenvector approximated the first order terms of the genie scores. This suggests that even when the likelihood is 801 not a linear function of the entries of the observation matrix, while the spectral algorithm may be 802 suboptimal, it could be still possible to get a non-trivial performance guarantee. 803

804 805

772

773

774

775

776

777

778

779 780

781

A.2 ADDITIONAL RELATED WORK

Community recovery under side information. Community detection problems with side information have been studied in numerous settings. Saad and Nosratinia Saad & Nosratinia (2018)
 considered exact recovery in the symmetric, balanced SBM, under the BEC, BSC, and more general side information with K features, where K may grow with n. Additionally, an efficient two-stage exact recovery algorithm was proposed. Vector-valued side information was also studied in Saad

810 & Nosratinia (2020), in the recovery of a planted dense subgraph of size o(n). Community detec-811 tion in the sparse setting under side information has received significant attention- see for example 812 Mossel & Xu (2016); Cai et al. (2016); Kadavankandy et al. (2017); Kanade et al. (2016); Desh-813 pande et al. (2018); Braun et al. (2022); we note that Deshpande et al. (2018) considers Gaussian 814 side information with either Bernoulli or Gaussian pairwise observations. See also Zhang et al. (2014) which includes statistical physics conjectures for recovery thresholds derived from the cavity 815 method. Numerous approaches for clustering have been proposed in the network science literature, 816 such as Newman & Clauset (2016); Zanghi et al. (2010); Yang et al. (2009); Xu et al. (2012); Yang 817 et al. (2013); Zhang et al. (2016); Gibert et al. (2012); Zhou et al. (2009); Günnemann et al. (2013); 818 Cheng et al. (2011); Binkiewicz et al. (2017); see Bothorel et al. (2015) for a survey. 819

820

842 843 844

845

821 Other inference problems with side information Related problems in the literature include doc-822 ument classification Krithara et al. (2008); Chang & Blei (2010) and text classification Balasubra-823 manyan & Cohen (2011). A recent line of work studies the problem of community detection from 824 correlated graphs Racz & Sridhar (2021); Gaudio et al. (2022), so that the additional graph plays the 825 role of side information. See also Zhang et al. (2021), which considers attributed graph alignment. More broadly, inference with side information falls under the area of semi-supervised learning (see 826 e.g. Van Engelen & Hoos (2020); Chapelle et al. (2002); Bair (2013); Basu et al. (2004); Newman 827 & Clauset (2016)). 828

Remark A.1 (\mathbb{Z}_2 -synchronization as rescaled ROS). We remark that the \mathbb{Z}_2 -synchronization problem is typically formulated as $A_{ij} = x_i^* x_j^* + \sigma W_{ij}$, where x^* is an unknown vector is chosen uniformly at random from the set $\{\pm 1\}^n$, W is a zero-diagonal symmetric matrix with independent entries sampled from $\mathcal{N}(0, 1)$ (Bandeira et al., 2017). In that case, the relevant parameterization of σ is $\sigma = c \sqrt{\frac{n}{\log n}}$, as $\sigma = \sqrt{\frac{n}{2\log n}}$ is the threshold value for exact recovery Bandeira et al. (2017). Thus, taking $a = 1/\sqrt{c}$, $b = -1/\sqrt{c}$ and $\rho = 1/2$ in our ROS model (Definition 2.2) produces a matrix A such that

$$A_{ij} = \frac{1}{c} \sqrt{\frac{\log n}{n}} x_i^* x_j^* + W, \quad x^* \sim \text{Uniform}(\{\pm 1\}^n).$$

After scaling A by $c\sqrt{n/\log n}$, we achieve the standard model $x^*x^{*\top} + \sigma W$ (with zero diagonal).

A.3 PROOF REVIEW OF DREVETON ET AL. (2024) FOR PROPOSITION 3.1

We first argue that the impossibility result in Proposition 3.1 is just a special case of the corresponding result in (Dreveton et al., 2024)[Theorem 1]. Note that we have only one pair of communities (C_+, C_-) , which is responsible for determining the threshold. Exactly as in Dreveton et al. (2024), we require the limit $L(t) = \lim_{n\to\infty} \frac{n}{\log n} \operatorname{CH}_t(+, -)$ to exist and be strictly concave.

For the MAP's success, we note that the only change is that we relax the requirement in two ways: (i) we do not require L(t) is strictly concave, (ii) we do not require that L(t) always exists, and are fine when the limit does not exist while diverging to $+\infty$. Dreveton et al. prove the MAP success in two parts, by first bounding the probability of the event that an assignment vector σ with hamming distance of m has higher likelihood than the ground truth σ^* (see their Lemma 3). In the second step (see their Appendix A.4), they do a union bound argument over all such possible vectors with the Hamming distance of m, and all varying $1 \le m \le n/2$.

We first note that their Lemma 3 statement as well as its proof do not make use of the condition that L(t) is strictly concave. In Appendix A.4, the union bound argument also does not require L(t)is strictly concave (it only uses $\sup_{t \in (0,1)} L(t) > 1$), and the strict concavity was only required to show the impossibility direction. To justify the second relaxation where we allow the cases of infinite limit L(t), we first note that their Lemma 3 is non-asymptotic anyway and derived for any finite *n*. Later, in Appendix A.4, we observe that the existence of the constants $\epsilon > 0$ and $\kappa > 0$ holds true, even when the limit $L(t) = +\infty$. In particular, the purpose behind the requirement was to exclude the cases when L(t) does not converge but neither approaches $+\infty$.

864 NOTATION AND PROBABILITY FACTS В 865

866 **B.1** NOTATION

867 Throughout the paper, we extensively use the standard extended real line algebra on $\mathbb{R} := \mathbb{R} \cup$ 868 $\{-\infty, +\infty\}$, where $+\infty$ and $-\infty$ are respectively greater and less than any other real number. Moreover, for any two $a, b \in \mathbb{R}$, we will define a - b = 0 when a = b. We also view $\log : \mathbb{R}_{>0} \to \mathbb{R}$ 870 where $\log(0) = -\infty$ and $\log(+\infty) = +\infty$, and use $0 \log 0 = 0 \log(+\infty) = 0$, following the 871 convention in the information-theory literature.

872 For any two vectors of random variables X and Y, we write $X \perp Y$ if entries of one are independent 873 from the others. For any real numbers $a, b \in R$, we denote $a \lor b = \max\{a, b\}$ and $a \land b = \min\{a, b\}$. 874 Let sgn : $\mathbb{R} \to \{\pm 1\}$ be the function defined by sgn(x) = 1 if $x \ge 0$ and sgn(x) = -1 if x < 0. 875 We also extend the definition to vectors; let sgn : $\mathbb{R}^n \to \{\pm 1\}^n$ be the map defined by applying the 876 sign function componentwise. We define $\mathbb{R}_+ = [0, \infty)$. For $n \in \mathbb{N}$, we write $[n] = \{1, 2, \dots, n\}$. We use the standard notation $o(.), O(.), \omega(.), \Omega(.), \Theta(.)$ etc. throughout the paper. For non-negative 877 sequences $(a_n)_{n\geq 1}$ and $(b_n)_{n\geq 1}$, we write $a_n \leq b_n$ to mean $a_n \leq Cb_n$ for some constant $\overline{C} > 0$. 878 The notation \asymp is similar, hiding two constants in upper and lower bounds. Moreover, we denote 879 $a_n \approx b_n$ as a shorthand for $\lim_{n\to\infty} \frac{a_n}{b_n} = 1$. 880

For a vector $x \in \mathbb{R}^n$, we define $||x||_2 = (\sum_{i=1}^n x_i^2)^{1/2}$, $||x||_1 = \sum_{i=1}^n |x_i|$, and $||x||_{\infty} = \max_i |x_i|$. Additionally, for any $i \in [n]$, we define x_{-i} as the vector in \mathbb{R}^n such that $(x_{-i})_j = x_j$ for $j \neq i$ 882 and $(x_{-i})_i = 0$. For any matrix $M \in \mathbb{R}^{n \times n}$, M_i refers to its ith row, which is a row vector, 883 and M_i refers to its ith column, which is a column vector. The matrix spectral norm is $||M||_2 =$ 884 $\sup_{\|x\|_{2}=1} \|Mx\|_{2}, \text{ the matrix } 2 \to \infty \text{ norm is } \|M\|_{2\to\infty} = \sup_{\|x\|_{2}=1} \|Mx\|_{\infty} = \sup_{i} \|M_{i}^{-}\|_{2}.$ 885 Let $zd : \mathbb{R}^{n \times n} \to \mathbb{R}^{n \times n}$ be the zero-diagonal mapping, where for any $A \in \mathbb{R}^{n \times n}$, $zd(A)_{ij} = A_{ij}$ if 886 $i \neq j$ and 0 otherwise. 887

B.2 STANDARD PROBABILITY LEMMAS

890 **Lemma B.1** (Rényi Divergence Formula). For any two probability laws $(\mathcal{A}, \mathcal{B})$, either both discrete 891 or continuous, we have that

$$e^{(t-1)D_t(\mathcal{A}||\mathcal{B})} = \mathbb{E}_{x \sim \mathcal{B}}\left[\left(\frac{\mathcal{A}(x)}{\mathcal{B}(x)}\right)^t\right]$$

Proof. By rearrangement the terms from the definition in (3).

Lemma B.2. Let $Z \sim \mathcal{N}(0, 1)$. The for any t > 0, 897

$$\left(\frac{1}{t} - \frac{1}{t^3}\right) \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \le \mathbb{P}(Z \ge t) \le \frac{1}{t\sqrt{2\pi}} e^{-t^2/2}.$$

We next show that the sampling procedure of the community assignment vector σ^* in any of the models leads to communities with roughly ρ and $(1 - \rho)$ fraction of vertices. 902

Lemma B.3. Let $\sigma^* \in \{\pm 1\}^n$ be a vector whose coordinates are i.i.d. with $\mathbb{P}(\sigma_i^* = +1) = \rho$. Then, let $C_+ := \{i : \sigma_i^* = +1\}$ and $C_- := \{i : \sigma_i^* = -1\}$. Define the event E as follow.

$$E := \{ ||C_+| - \rho n| \le \rho n^{2/3} \quad \text{and} \quad ||C_-| - (1 - \rho)n| \le \rho n^{2/3} \}.$$

$$E) \ge 1 - o(1).$$
(9)

Then $\mathbb{P}(E) \geq$ 906

888

889

896

898 899 900

901

903

904

905

911

913 914

907 *Proof.* For each $1 \leq i \leq n$, since $\mathbb{P}(\sigma_i^*) = \rho$ i.i.d., we have that $|C_+| = |\{i: \sigma_i^* = +1\}|$ fol-908 lows $Bin(n, \rho)$. The Chernoff bound for binomial random variables (Mitzenmacher & Upfal, 2017, 909 Theorem 4.4, Theorem 4.5) implies that for any $\delta \in (0, 1)$, 910

$$\mathbb{P}(||C_+| - \mathbb{E}[|C_+|]| \le \delta n) \ge 1 - 2\exp\left(-\delta^2 \mathbb{E}[|C_+|]/3\right)$$

Note that $\mathbb{E}[|C_+|] = \rho n$. Choosing $\delta = \rho n^{-1/3} = o(1)$ (as $\rho > 0$ is a constant) 912

$$\mathbb{P}(||C_+| - \rho n| \le \rho n^{2/3}) \ge 1 - 2\exp\left(-n^{-2/3}\rho^3 n/3\right) = 1 - 2\exp\left(-\rho^3 n^{1/3}/3\right).$$

Since $n^{1/3} = \omega(\log n)$ and $\rho > 0$ is a constant, 915

916
917
$$\mathbb{P}\left((1-n^{-1/3})\,\rho n \le |C_+| \le \left(1+n^{-1/3}\right)\rho n\right) \ge 1 - O(\exp\left(-10\log n/3\right)) = 1 - O(n^{-3}),$$
which implies the lemma.

17

918 C VERIFICATION FOR INFORMATION THEORETIC LIMITS

920 In this section, we will verify the function L(t) is well defined and is strictly concave for both ROS 921 and SBM and in the cases of no side information, or the special cases of GF, BEC and BSC side 922 information channels. As a result, even I^* is well-defined for these settings and sharply characterizes 923 the information-theoretic limits in these settings, making our spectral algorithms in Theorem 1 and 924 Theorem 2 optimal, when combined with the impossibility direction of Proposition 3.1. In order 925 to verify these technical conditions, we will use a well-established closed form expression of Rényi 926 divergence of two Gaussian random variables, e.g. see Gil et al. (2013).

1927 Lemma C.1. Let $\mathcal{A} \equiv \mathcal{N}(\mu_1, \sigma^2 I_d)$ and $\mathcal{B} \equiv \mathcal{N}(\mu_2, \sigma^2 I_d)$ for $\mu_1, \mu_2 \in \mathbb{R}^d$ and $\sigma^2 > 0$, then

$$D_t(\mathcal{A} \| \mathcal{B}) = rac{t}{2\sigma^2} \| \mu_1 - \mu_2 \|_2^2.$$

931 Side Information Terms. We start by analyzing the additive terms that come in L(t) due to the 932 presence of each example of side information channels.

Lemma C.2. Consider the side information laws $S_+ \equiv \mathcal{N}(v_+, \sigma^2 I_d)$ and $S_- \equiv \mathcal{N}(v_-, \sigma^2 I_d)$, for constant $\sigma^2 > 0$ and $v_+ = \beta_+ \sqrt{\log n}$ and $v_- = \beta_- \sqrt{\log n}$ for $\beta_+, \beta_- \in \mathbb{R}^d$. Then

$$\lim_{n \to \infty} \frac{n}{\log n} \cdot \frac{(1-t)}{n} D_t(\mathcal{S}_+ \| \mathcal{S}_-) = \frac{t(1-t) \left\| \mathcal{\beta}_+ - \mathcal{\beta}_- \right\|_2^2}{2\sigma^2}$$

which is well-defined. Note that this as a function of $t \in (0,1)$ is strictly concave when $\beta_+ \neq \beta_-$ or it is 0 otherwise.

Proof. The lemma follows from a straighforward simplification and the use of Lemma C.1

$$\lim_{n \to \infty} \frac{n}{\log n} \cdot \frac{(1-t)}{n} \mathcal{D}_t(\mathcal{S}_+ \| \mathcal{S}_-) = \lim_{n \to \infty} \frac{(1-t)}{\log n} \frac{t \| v_+ - v_- \|_2^2}{2\sigma^2}$$
$$= \lim_{n \to \infty} \frac{t(1-t)}{\log n} \frac{\| \beta_+ - \beta_- \|_2^2 \log n}{2\sigma^2}$$
$$= \frac{t(1-t) \| \beta_+ - \beta_- \|_2^2}{2\sigma^2}.$$

Lemma C.3. Consider the side information laws (S_+, S_-) as in the BEC channel with ϵ such that $\lim_{n\to\infty} \frac{\log(1/\epsilon)}{\log n} = \beta$ for $\beta \ge 0$, i.e. $\epsilon = n^{-\beta+o(1)}$ Then

$$\lim_{n \to \infty} \frac{n}{\log n} \times \frac{(1-t)}{n} D_t(\mathcal{S}_+ \| \mathcal{S}_-) = \beta$$

Proof. By definition of Rényi divergence in (3)

$$\lim_{n \to \infty} \frac{n}{\log n} \cdot \frac{(1-t)}{n} \mathcal{D}_t(\mathcal{S}_+ \| \mathcal{S}_-) = \lim_{n \to \infty} \frac{(1-t)}{\log n} \frac{1}{(t-1)} \log_{y_i \sim \mathcal{S}_-} \left[\left(\frac{\mathcal{S}_+(y_i)}{\mathcal{S}_-(y_i)} \right)^t \right]$$
$$= \lim_{n \to \infty} -\frac{\log \epsilon}{\log n} = \lim_{n \to \infty} \frac{\log(1/\epsilon)}{\log n} = \beta.$$

Lemma C.4. Consider the side information laws (S_+, S_-) as in the BSC channel with α such that $\lim_{n\to\infty} \frac{\log(\frac{1-\alpha}{\alpha})}{\log n} = \beta$ for $\beta \ge 0$, i.e. $\alpha = n^{-\beta+o(1)}$ Then

$$\lim_{n \to \infty} \frac{n}{\log n} \times \frac{(1-t)}{n} D_t(\mathcal{S}_+ \| \mathcal{S}_-) = \beta \min\{t, 1-t\}$$

Note that this is a concave function of $t \in (0, 1)$.

Proof. By definition of Rényi divergence in (3)

$$\lim_{n \to \infty} \frac{n}{\log n} \cdot \frac{(1-t)}{n} \mathcal{D}_t(\mathcal{S}_+ \| \mathcal{S}_-) = \lim_{n \to \infty} \frac{(1-t)}{\log n} \frac{1}{(t-1)} \log \mathbb{E}_{y_i \sim \mathcal{S}_-} \left[\left(\frac{\mathcal{S}_+(y_i)}{\mathcal{S}_-(y_i)} \right)^t \right]$$
$$= \lim_{n \to \infty} -\frac{\log(\alpha^t (1-\alpha)^{1-t} + \alpha^{1-t} (1-\alpha)^t)}{\log n}$$
$$= -\lim_{n \to \infty} \frac{-\min\{\beta t, \beta (1-t)\} \log n}{\log n} = \beta \min\{t, 1-t\}.$$

> **Deriving Threshold for** ROS and SBM Finally, we will derive the expression of L(t) and show that it is strictly concave.

> **Lemma C.5.** Fix $\rho \in (0,1)$ and $a, b \in \mathbb{R}$ such that $\max\{|a|, |b|\} > 0$. Consider the distribution $(\mathcal{P}_+, \mathcal{P}_-, \mathcal{Q})$ as defined in the ROS model (Definition 2.2). Then in the absence of side information, L(t) is well-defined and given by

$$L(t) = t(1-t)\frac{(a-b)^2(\rho a^2 + (1-\rho)b^2)}{2}$$

Moreover, L(t) is strictly concave. Additionally, under GF, BEC, BSC channels described in Lemmas C.2, C.3, and C.4 respectively, L(t) continues to be well-defined and strictly concave.

Proof. In the absence of side information

$$L(t) = \lim_{n \to \infty} \frac{n}{\log n} \operatorname{CH}_t(+, -) = \lim_{n \to \infty} \frac{n(1-t)}{\log n} \left[\rho \operatorname{D}_t(\mathcal{P}_+ \| \mathcal{Q}) + (1-\rho) \operatorname{D}_t(\mathcal{Q} \| \mathcal{P}_-) \right]$$

$$= \lim_{n \to \infty} \frac{n(1-t)}{\log n} \left[\rho \cdot \frac{t(a^2 - ab)^2}{2} \frac{\log n}{n} + (1-\rho) \cdot \frac{t(ab - b^2)^2}{2} \frac{\log n}{n} \right]$$

(using Lemma B.1)

$$= t(1-t)\frac{\left(\rho(a^2-ab)^2 + (1-\rho)(ab-b^2)^2\right)}{2} = t(1-t)\frac{(a-b)^2(\rho a^2 + (1-\rho)b^2)}{2}.$$

Note that when $\max\{|a|, |b|\} > 0$, the multiplicative coefficient of t(1-t) is positive. Therefore, L(t) is strictly concave function in (0, 1).

Under side information, L(t) is given by adding another term based on the channel as derived in Lemmas C.2, C.3, and C.4. Therefore, L(t) continues to be well-defined. Moreover, L(t) still remains strictly concave because the addition of a strictly concave function with another concave function (including just constant) is strictly concave. П

Lemma C.6. Fix $\rho \in (0,1)$ and $a_1, a_2, b > 0$. Consider the distribution $(\mathcal{P}_+, \mathcal{P}_-, \mathcal{Q})$ as defined in the SBM model (Definition 2.3). Then in the absence of side information, L(t) is well-defined and given by

$$L(t) = t\rho a_1 + (1-t)\rho b + t(1-\rho)b + (1-t)(1-\rho)a_2 - \rho a_1^t b^{1-t} - (1-\rho)b^t a_2^{1-t}.$$

Moreover, except the degenerate case $a_1 = a_2 = b$, we have that L(t) is strictly concave. Additionally, under GF, BEC, BSC channels described in Lemmas C.2, C.3, and C.4 respectively, L(t)continues to be well-defined and strictly concave.

Proof. This is just a special case of the expression derived by (Dreveton et al., 2024, Example 1) but for two communities. Note that L(t) is twice differentiable, and

$$L''(t) = -\rho\left(\frac{a_1}{b}\right)^t b \log^2\left(\frac{a_1}{b}\right) - (1-\rho)\left(\frac{b}{a_2}\right)^t a_2 \log^2\left(\frac{b}{a_2}\right) < 0$$

unless $a_1 = a_2 = b$, and thus, L(t) is strictly concave. Similar to Lemma C.5, even for SBM, after adding another side information channel-specific term to L(t), derived in Lemmas C.2, C.3, and C.4, the function L(t) is well-defined and strictly concave.

1026 D **OMITTED PROOFS FROM SECTION 4** 1027

1028 In this section, we prove all the genie-aided estimation related lemmas. We start with the proof of 1029 Lemma 4.1, which establishes that the failure of some genie-aided estimator implies that the global 1030 MAP estimator also fails.

1031 1032

1039 1040

104

105

1062

1063

1065

1066

1067

1068

1033 *Proof of Lemma 4.1.* We first consider the case when side information y (either BEC or BSC) is 1034 provided. Let $S_1 = \{(A, \overline{\sigma}, \overline{y}) : \exists i \in [n] \text{ such that } \hat{\sigma}_{\text{Gen}, i} \text{ fails}\}$. Similarly, $S_2 = \{(A, \overline{\sigma}, \overline{y}) : \exists i \in [n] \text{ such that } \hat{\sigma}_{\text{Gen}, i} \text{ fails}\}$. 1035 $\hat{\sigma}_{MAP}$ fails}. To show the desired claim, it suffices to show that $S_1 \subseteq S_2$. 1036

To this end, fix any instance $(\overline{A}, \overline{\sigma}, \overline{y}) \in S_1$ of (A, σ^*, y) . Then by definition of the failure of the 1037 genie-aided estimators give below Equation (6), there exists $i \in [n]$ such that 1038

$$\mathbb{P}(\sigma_i^* = -\overline{\sigma}_i \mid \overline{A}, \overline{y}, \overline{\sigma}_{-i}) > \mathbb{P}(\sigma_i^* = \overline{\sigma}_i \mid \overline{A}, \overline{y}, \overline{\sigma}_{-i}).$$
(10)

1041 We now consider the community assignment vector $\sigma' \in \{\pm 1\}^n$ whose labeling agrees with $\overline{\sigma}$ 1042 except for the i^{th} label, for which $\sigma'_i = -\overline{\sigma}_i$. Then 1043

$$\begin{array}{ll} 1044 \\ 1045 \\ 1046 \\ 1046 \\ 1046 \\ 1046 \\ 1046 \\ 1046 \\ 1046 \\ 1047 \\ 1047 \\ 1047 \\ 1047 \\ 1048 \\ 1049 \\ 1050 \\ 1050 \\ 1050 \\ 1051 \\ 1052 \\ 1052 \\ 1053 \end{array} \\ \begin{array}{ll} \mathbb{P}(\sigma^*_{-i} = \sigma'_{-i} \mid \overline{A}, \overline{y}) \cdot \mathbb{P}(\sigma^*_i = \sigma'_i \mid \overline{A}, \overline{y}, \sigma^*_{-i} = \overline{\sigma}_{-i}) \\ \mathbb{P}(\sigma^*_{-i} = \overline{\sigma}_{-i} \mid \overline{A}, \overline{y}) \cdot \mathbb{P}(\sigma^*_i = -\overline{\sigma}_i \mid \overline{A}, \overline{y}, \sigma^*_{-i} = \overline{\sigma}_{-i}) \\ \mathbb{P}(\sigma^*_{-i} = \overline{\sigma}_{-i} \mid A, y) \cdot \mathbb{P}(\sigma^* = \overline{\sigma}_i \mid \overline{A}, \overline{y}, \sigma^*_{-i} = \overline{\sigma}_{-i}) \\ \mathbb{P}(\sigma^*_{-i} = \overline{\sigma}_{-i}, \sigma^*_i = \overline{\sigma}_i \mid \overline{A}, \overline{y}) \\ \mathbb{P}(\sigma^*_{-i} = \overline{\sigma}_{-i}, \sigma^*_i = \overline{\sigma}_i \mid \overline{A}, \overline{y}) \\ \mathbb{P}(\sigma^*_{-i} = \overline{\sigma}_{-i}, \sigma^*_i = \overline{\sigma}_i \mid \overline{A}, \overline{y}) \\ \mathbb{P}(\sigma^*_{-i} = \overline{\sigma}_{-i}, \sigma^*_i = \overline{\sigma}_i \mid \overline{A}, \overline{y}) \\ \mathbb{P}(\sigma^*_{-i} = \overline{\sigma}_{-i} \mid \overline{A}, \overline{y}) . \end{array}$$

1054 By the definition of the MAP estimator $\hat{\sigma}_{MAP} = \arg \max_{\sigma \in \{\pm 1\}^n} \mathbb{P}(\sigma^* = \sigma | \overline{A}, \overline{y}) \neq \overline{\sigma}$, which 1055 implies $(A, \overline{\sigma}, \overline{y}) \in \mathcal{S}_2$. 1056

Finally, we consider the case when there is no side information. Define S_1 and S_2 similarly, but 1057 after dropping y. 1058

- Case $\mathcal{P}_+ \neq \mathcal{P}_-$ or $\rho \neq 1/2$. Consider any $(\overline{A}, \overline{\sigma}) \in \mathcal{S}_1$ and follow exactly the same argument used in in deriving (11) (after dropping the conditioning on y). This will lead to the conclusion that $(A, \overline{\sigma}) \in S_2$, yielding $S_1 \subseteq S_2$.
- Case $\mathcal{P}_+ \equiv \mathcal{P}_-$ and $\rho = 1/2$. Consider any $(\overline{A}, \overline{\sigma}) \in \mathcal{S}_1$. Follow exactly the same argument as in (11) and conclude that $\mathbb{P}(\sigma^* = \sigma' \mid \overline{A}) > \mathbb{P}(\sigma^* = \overline{\sigma} \mid \overline{A})$. Additionally, due to the symmetry, $\mathbb{P}(\sigma^* = \overline{\sigma} \mid \overline{A}) = \mathbb{P}(\sigma^* = -\overline{\sigma} \mid A)$. Combining these two, we obtain that $\hat{\sigma}_{MAP} \notin \{\pm \sigma^*\}$, which implies $\hat{\sigma}_{MAP}$ fails by Definition 2.8. Conclude $(\overline{A}, \overline{\sigma}) \in S_2$, as desired.

We now derive the expressions for genie scores given by Lemma 4.2, without or with side informa-1075 tion. 1076

1077 1078

1074

Proof of Lemma 4.2. We first recall the definition of genie scores from (7). For any $i \in [n]$, we do 1079 the following analysis in each model of side information.

 $= \log \left(\frac{\mathbb{P}(\sigma_i^* = +1) \cdot \mathcal{L}(A \mid \sigma_i^* = +1, \sigma_{-i}^*)}{\mathbb{P}(\sigma_i^* = -1) \cdot \mathcal{L}(A \mid \sigma_i^* = -1, \sigma_{-i}^*)} \right)$

No side information:

 $= \log \left(\frac{\mathbb{P}(\sigma_i^* = +1) \cdot \mathcal{L}(A_i. \mid \sigma_i^* = +1, \sigma_{-i}^*)}{\mathbb{P}(\sigma_i^* = -1) \cdot \mathcal{L}(A_i. \mid \sigma_i^* = -1, \sigma_{-i}^*)} \right)$ (the likelihood of all but the *i*th row is same conditioned under σ_{-i}^* irrespective of σ_i^*)

$$= \log \left(\rho \cdot \prod_{j \in C_{+}^{-i}} \mathcal{P}_{+}(A_{ij}) \cdot \prod_{j \in C_{-}^{-i}} \mathcal{Q}(A_{ij}) \middle/ (1-\rho) \cdot \prod_{j \in C_{+}^{-i}} \mathcal{Q}(A_{ij}) \cdot \prod_{j \in C_{-}^{-i}} \mathcal{P}_{-}(A_{ij}) \right)$$

(due to the conditional independence of the entries and the law of GBM)

$$= \sum_{j \in C_+^{-i}} \log\left(\frac{\mathcal{P}_+(A_{ij})}{\mathcal{Q}(A_{ij})}\right) + \sum_{j \in C_-^{-i}} \log\left(\frac{\mathcal{Q}(A_{ij})}{\mathcal{P}_-(A_{ij})}\right) + \log\left(\frac{\rho}{1-\rho}\right).$$

 $z_i^* = \log\left(\frac{\mathbb{P}(\sigma_i^* = +1 \mid A, \sigma_{-i}^*)}{\mathbb{P}(\sigma_i^* = -1 \mid A, \sigma_{-i}^*)}\right) = \log\left(\frac{\mathbb{P}(\sigma_i^* = +1) \cdot \mathcal{L}(A, \sigma_{-i}^* \mid \sigma_i^* = +1)}{\mathbb{P}(\sigma_i^* = -1) \cdot \mathcal{L}(A, \sigma_{-i}^* \mid \sigma_i^* = -1)}\right)$

Side information: Again by (7), we have

$$z_i^* = \log \left(\frac{\mathbb{P}(\sigma_i^* = +1 \mid A, y, \sigma_{-i}^*)}{\mathbb{P}(\sigma_i^* = -1 \mid A, y, \sigma_{-i}^*)} \right) = \log \left(\frac{\mathbb{P}(\sigma_i^* = +1 \mid A, y_i, \sigma_{-i}^*)}{\mathbb{P}(\sigma_i^* = -1 \mid A, y_i, \sigma_{-i}^*)} \right)$$
(conditioned on σ_{-i}^* , we have $\sigma_i^* \perp y_{-i}$)

$$= \log \left(\frac{\mathbb{P}(\sigma_i^* = +1) \cdot \mathcal{L}(A, y_i, \sigma_{-i}^* \mid \sigma_i^* = +1)}{\mathbb{P}(\sigma_i^* = -1) \cdot \mathcal{L}(A, y_i, \sigma_{-i}^* \mid \sigma_i^* = -1)} \right)$$

$$= \log \left(\frac{\mathbb{P}(\sigma_i^* = +1) \cdot \mathcal{L}(A, \sigma_{-i}^* \mid \sigma_i^* = +1) \cdot \mathbb{P}(y_i \mid \sigma_i^* = +1)}{\mathbb{P}(\sigma_i^* = -1) \cdot \mathcal{L}(A, \sigma_{-i}^* \mid \sigma_i^* = -1) \cdot \mathbb{P}(y_i \mid \sigma_i^* = -1)} \right)$$

(conditioned on σ_i^* , we have $y_i \perp A$ and $y_i \perp \sigma_{-i}^*$)

 $(\sigma_{-i}^* \text{ is independent of } \sigma_i^*)$

$$= \log \left(\frac{\mathbb{P}(\sigma_i^* = +1 \mid A, \sigma_{-i}^*) \cdot \mathbb{P}(y_i \mid \sigma_i^* = +1)}{\mathbb{P}(\sigma_i^* = -1 \mid A, \sigma_{-i}^*) \cdot \mathbb{P}(y_i \mid \sigma_i^* = -1)} \right)$$

$$(\mathbb{P}(\sigma_i^* = +1 \mid A, \sigma_{-i}^*)) = (\mathbb{P}(\alpha_i \mid \sigma_i^* = -1))$$

$$= \log\left(\frac{\mathbb{P}(\sigma_{i}^{*} = +1 \mid A, \sigma_{-i}^{*})}{\mathbb{P}(\sigma_{i}^{*} = -1 \mid A, \sigma_{-i}^{*})}\right) + \log\left(\frac{\mathbb{P}(y_{i} \mid \sigma_{i}^{*} = +1)}{\mathbb{P}(y_{i} \mid \sigma_{i}^{*} = -1)}\right)$$
(12)

Note that

$$\log\left(\frac{\mathbb{P}(y_i \mid \sigma_i^* = +1)}{\mathbb{P}(y_i \mid \sigma_i^* = -1)}\right) = \log\left(\frac{\mathcal{S}_+(y_i)}{\mathcal{S}_-(y_i)}\right).$$

Substituting this in (12) along with the definition of genie score without side information (7) and using the expression from the case without side information, we obtain

$$z_{i}^{*} = \sum_{j \in C_{+}^{-i}} \log\left(\frac{\mathcal{P}_{+}(A_{ij})}{\mathcal{Q}(A_{ij})}\right) + \sum_{j \in C_{-}^{-i}} \log\left(\frac{\mathcal{Q}(A_{ij})}{\mathcal{P}_{-}(A_{ij})}\right) + \log\left(\frac{\rho}{1-\rho}\right) + \log\left(\frac{\mathcal{S}_{+}(y_{i})}{\mathcal{S}_{-}(y_{i})}\right). \quad \Box$$

$$1122$$

We now show the proof of Lemma 4.3, where we explicitly derive the closed form of the addi-tive factor under side information for the special cases of Gaussian Features (GF), Binary Erasure Channel (BEC), and Binary Symmetric Channel (BSC).

Proof of Lemma 4.3. We break into cases.

• Gaussian Features (GF):

$$\log\left(\frac{\mathcal{S}_{+}(y_{i})}{\mathcal{S}_{-}(y_{i})}\right) = \log\left(\frac{e^{-\frac{\|y_{i}-\upsilon_{+}\|_{2}^{2}}{2\sigma^{2}}}}{e^{-\frac{\|y_{i}-\upsilon_{-}\|_{2}^{2}}{2\sigma^{2}}}}\right) = \log\left[\exp\left(\frac{\|y_{i}-\upsilon_{-}\|_{2}^{2}-\|y_{i}-\upsilon_{+}\|_{2}^{2}}{2\sigma^{2}}\right)\right]$$

1133
$$= \frac{\|y_i - v_-\|_2^2 - \|y_i - v_+\|_2^2}{2\sigma^2}.$$

...

• Binary Erasure Channel (BEC): If $y_i = 0$, then $S_+(y_i) = S_-(y_i) = \epsilon$. If $y_i = +1$, then $S_+(y_i) = 1 - \epsilon$ but $S_-(y_i) = 0$, and similarly, if $y_i = -1$, then $S_+(y_i) = 0$ but $S_-(y_i) = 1 - \epsilon$.

$$\log\left(\frac{\mathcal{S}_{+}(y_{i})}{\mathcal{S}_{-}(y_{i})}\right) = \begin{cases} +\infty, & \text{if } \sigma_{i}^{*} = +1; \\ -\infty, & \text{if } \sigma_{i}^{*} = -1; \\ 0, & \text{otherwise.} \end{cases}$$

• Binary Symmetric Channel (BSC):

$$\log\left(\frac{\mathcal{S}_{+}(y_{i})}{\mathcal{S}_{-}(y_{i})}\right) = \log\left(\frac{1-\alpha}{\alpha}\right)\mathbf{1}[y_{i}=+1] + \log\left(\frac{\alpha}{1-\alpha}\right)\mathbf{1}[y_{i}=-1] = \log\left(\frac{1-\alpha}{\alpha}\right)y_{i}.$$

1148 D.1 GENIE SUCCESS WITH MARGIN ABOVE THE IT THRESHOLD.

¹¹⁴⁹ We now give a formal version of the claim that above the information-theoretic limit, all the genie scores succeed with a margin of $\Omega(\log n)$, formalized in the following lemma.

Lemma D.1. Fix $\rho \in (0,1)$ and consider probability laws $(\mathcal{P}_+, \mathcal{P}_-, \mathcal{Q})$. Let $(A, \sigma^*) \sim GBM_n(\rho, \mathcal{P}_+, \mathcal{P}_-, \mathcal{Q})$ and we observe A. Condition on σ^* such that the event E from (9) holds. Optionally, for side information laws $(\mathcal{S}_+, \mathcal{S}_-)$, we observe $y \sim Sl(\sigma^*, \mathcal{S}_+, \mathcal{S}_-)$. Let z^* be the genie score vector for the corresponding model given by (7). Let I^* be defined according to (5). Then if $I^* > 1$ then there exists some constant $\delta > 0$ such that for any $i \in [n]$:

$$\mathbb{P}(\sigma_i^* z_i^* < \delta \log n) = o(n^{-1}).$$

Proof. By Lemma 4.2, for any $i \in [n]$, the genie score z_i^* is given by

$$z_i^* = \sum_{j \in C_+^{-i}} \log\left(\frac{\mathcal{P}_+(A_{ij})}{\mathcal{Q}(A_{ij})}\right) + \sum_{j \in C_-^{-i}} \log\left(\frac{\mathcal{Q}(A_{ij})}{\mathcal{P}_-(A_{ij})}\right) + \log\left(\frac{\rho}{1-\rho}\right) + \log\left(\frac{\mathcal{S}_+(y_i)}{\mathcal{S}_-(y_i)}\right)$$

For convenience, we define $X_i := z_i^* - \log\left(\frac{\rho}{1-\rho}\right)$. For any $i \in C_+$ and any $\varepsilon, t > 0$,

$$\begin{array}{ll} \mathbf{1164} & \mathbb{P}(\sigma_i^* z_i^* < \varepsilon \log n) = \mathbb{P}(z_i^* < \varepsilon \log n) = \mathbb{P}\left(X_i < (1 + o(1))\varepsilon \log n\right) = \mathbb{P}\left(e^{tX_i} < e^{t(1 + o(1))\varepsilon \log n}\right) \\ & \leq e^{t\varepsilon \log n} \mathbb{E}\left[e^{-tX_i}\right]. \end{array}$$

1167 We now analyze

$$\begin{aligned} & \begin{array}{l} \mathbf{1168} \\ \mathbf{1169} \\ \mathbf{1170} \\ \mathbf{1170} \\ \mathbf{1170} \\ \mathbf{1171} \\ \end{array} & \begin{array}{l} \mathbb{E}\left[e^{-tX_{i}}\right] = \mathbb{E}\left[\exp\left(-t\left(\sum_{j\in C_{+}^{-i}}\log\left(\frac{\mathcal{P}_{+}(A_{ij})}{\mathcal{Q}(A_{ij})}\right) + \sum_{j\in C_{-}^{-i}}\log\left(\frac{\mathcal{Q}(A_{ij})}{\mathcal{P}_{-}(A_{ij})}\right) + \log\left(\frac{\mathcal{S}_{+}(y_{i})}{\mathcal{S}_{-}(y_{i})}\right)\right)\right) \\ \mathbf{1171} \\ \mathbf{1172} \\ \mathbf{1173} \\ \mathbf{1174} \\ \mathbf{1175} \\ \end{array} & = \mathbb{E}\left[\exp\left(t\left(\sum_{j\in C_{+}^{-i}}\log\left(\frac{\mathcal{Q}(A_{ij})}{\mathcal{P}_{+}(A_{ij})}\right) + \sum_{j\in C_{-}^{-i}}\log\left(\frac{\mathcal{P}_{-}(A_{ij})}{\mathcal{Q}(A_{ij})}\right) + \log\left(\frac{\mathcal{S}_{-}(y_{i})}{\mathcal{S}_{+}(y_{i})}\right)\right)\right)\right) \\ \mathbf{1176} \\ \mathbf{1176} \\ \mathbf{1178} \\ \end{array} & = \mathbb{E}\left[\prod_{j\in C_{+}^{-i}}e^{t\log\left(\frac{\mathcal{Q}(A_{ij})}{\mathcal{P}_{+}(A_{ij})}\right)}\prod_{j\in C_{-}^{-i}}e^{t\log\left(\frac{\mathcal{P}_{-}(A_{ij})}{\mathcal{Q}_{+}(A_{ij})}\right)} \cdot e^{t\log\left(\frac{\mathcal{S}_{-}(y_{i})}{\mathcal{S}_{+}(y_{i})}\right)}\right] \\ \end{array}$$

$$= \prod_{j \in C_{+}^{-i}} \mathbb{E}\left[\left(\frac{\mathcal{Q}(A_{ij})}{\mathcal{P}_{+}(A_{ij})}\right)^{t}\right] \prod_{j \in C_{-}^{-i}} \mathbb{E}\left[\left(\frac{\mathcal{P}_{-}(A_{ij})}{\mathcal{Q}(A_{ij})}\right)^{t}\right] \mathbb{E}\left[\left(\frac{\mathcal{S}_{-}(y_{i})}{\mathcal{S}_{+}(y_{i})}\right)^{t}\right]$$

$$= e^{-(1+o(1))n\operatorname{CH}_{1-t}(+,-)}$$

1188 Substituting this in our Chernoff-style bound, we obtain for any $i \in C_+$ 1189 $\mathbb{P}(\sigma_i^* z_i^* < \varepsilon \log n) \le e^{-(1+o(1))(n\operatorname{CH}_{1-t}(+,-)-t\varepsilon \log n)}$ 1190 When $I^* > 1$, we have there exists $t^* \in (0, 1)$ such that $L(1-t^*) = \lim_{n \to \infty} \frac{n}{\log n} \operatorname{CH}_{1-t^*}(+, -) > 0$ 1191 1. Thus, there exists a constant $\epsilon > 0$, and $t^* \in (0, 1)$ such that for sufficient large n: 1192 1193 $\mathbb{P}(\sigma_i^* z_i^* < \varepsilon \log n) < e^{-((1+\epsilon)\log n - t^* \varepsilon \log n)}.$ 1194 Therefore, one can choose $\delta_1 := \delta_1(t^*, \epsilon) > 0$ small enough such that 1195 $\mathbb{P}(\sigma_i^* z_i^* < \delta_1 \log n) \le e^{-((1+\epsilon/2)\log n)} = o(n^{-1}).$ 1196 We carry out exactly similar calculation for the community C_- . For any $i \in C_-$ and $t, \varepsilon > 0$, 1197 1198 $\mathbb{P}(\sigma_i^* z_i^* < \varepsilon \log n) = \mathbb{P}(z_i^* > -\varepsilon \log n) = \mathbb{P}\left(X_i > -(1+o(1))\varepsilon \log n\right) = \mathbb{P}\left(e^{tX_i} > e^{-t(1+o(1))\varepsilon \log n}\right)$ 1199 $< e^{t\varepsilon \log n} \mathbb{E}\left[e^{tX_i}\right].$ 1200 1201 Simplifying 1202 $\mathbb{E}\left[e^{tX_i}\right] = \mathbb{E}\left|\exp\left(t\left(\sum_{j\in C_{\perp}^{-i}}\log\left(\frac{\mathcal{P}_+(A_{ij})}{\mathcal{Q}(A_{ij})}\right) + \sum_{j\in C_{\perp}^{-i}}\log\left(\frac{\mathcal{Q}(A_{ij})}{\mathcal{P}_-(A_{ij})}\right) + \log\left(\frac{\mathcal{S}_+(y_i)}{\mathcal{S}_-(y_i)}\right)\right)\right)\right|$ 1203 1204 1205 1206 $=\prod_{j\in C_{-}^{-i}} \mathbb{E}\left[\left(\frac{\mathcal{P}_{+}(A_{ij})}{\mathcal{Q}(A_{ij})}\right)^{t}\right] \prod_{j\in C_{-}^{-i}} \mathbb{E}\left[\left(\frac{\mathcal{Q}(A_{ij})}{\mathcal{P}_{-}(A_{ij})}\right)^{t}\right] \mathbb{E}\left[\left(\frac{\mathcal{S}_{+}(y_{i})}{\mathcal{S}_{-}(y_{i})}\right)^{t}\right]$ 1207 1208 1209 $- \rho^{(1+o(1))(t-1)(\rho n \mathsf{D}_t(\mathcal{P}_+ \| \mathcal{Q}) + (1-\rho)n \mathsf{D}_t(\mathcal{Q} \| \mathcal{P}_-) + \mathsf{D}_t(\mathcal{S}_+ \| \mathcal{S}_-))}$ 1210 (using Lemma B.1 and using community sizes conditioned on E) 1211 $= e^{-(1+o(1))nCH_t(+,-)}$ 1212 1213 Overall, we obtain for any $i \in C_+$ 1214 $\mathbb{P}(\sigma_i^* z_i^* < \varepsilon \log n) \le e^{-(1+o(1))(n \operatorname{CH}_t(+,-) - t\varepsilon \log n)}$ 1215 If $I^* > 1$, then there exists $t^* \in (0,1)$ such that $L(t^*) = \lim_{n \to \infty} \frac{n}{\log n} \operatorname{CH}_t(+,-) > 1$ and thus, 1216 there exists a constant $\epsilon > 0$, and $t^* \in (0, 1)$ such that for sufficient large n: 1217 1218 $\mathbb{P}(\sigma_i^* z_i^* < \varepsilon \log n) \le e^{-((1+\epsilon)\log n - t^* \varepsilon \log n)}$ 1219 Therefore, choosing $\delta_2 := \delta_2(t^*, \epsilon) > 0$ small enough such that 1220 $\mathbb{P}(\sigma_i^* z_i^* < \delta_2 \log n) \le e^{-((1+\epsilon/2)\log n)} = o(n^{-1}).$ 1221 1222 Finally, choosing $\delta = \min{\{\delta_1, \delta_2\}} > 0$, we obtain for $i \in [n]$ 1223 $\mathbb{P}(\sigma_i^* z_i^* < \delta \log n) = o(n^{-1}),$ 1224 concluding the proof of the lemma. 1225 1226 ENTRYWISE BEHAVIOR OF EIGENVECTORS. Е 1227 1228 Abbe, Fan, Wang, and Zhong Abbe et al. (2020) showed the powerful entrywise behavior of eigenvectors for a general ensemble of random matrices under certain assumptions. Their result (Abbe 1230 et al., 2020, Theorem 2.1) applies more generally to eigenspaces; below we note a special case of

1231 their result when the eigenspace has a single eigenvalue. 1232 Suppose $A \in \mathbb{R}^{n \times n}$ is a symmetric random matrix and $A^* = \mathbb{E}[A]$. Let the eigenvalues of A1233 be $|\lambda_1| \ge \cdots \ge |\lambda_n|$, and their associated eigenvectors be $\{u_j\}_{j\in[n]}$ (defined up to rotation if 1234 eigenvalues are repeated). Analogously for A^* , the eigenvalues and eigenvectors are $|\lambda_1^*| \ge \cdots \ge$ 1235 $|\lambda_n^*|$ and $\{u_i^*\}_{j\in[n]}$, respectively. For any fixed (λ_i^*, u_i^*) , define the eigengap quantity

1236 1237

$$\Delta^* := |\lambda_i^*| \wedge \min_{j \in [n] \setminus \{i\}} |\lambda_i^* - \lambda_j^*|.$$
(13)

Here we define the eigengap for the special case of (Abbe et al., 2020, Theorem 2.1) applied to a single eigenvector, rather than for an eigenspace associated with consecutive eigenvalues. For more general definition when the eigenspace contains multiple eigenvalues, see (Abbe et al., 2020, Equation (2.1)). We define $\kappa := |\lambda_i^*|/\Delta^*$, which is always bounded from below by 1. For a parameter $\gamma \ge 0$, consider the following four assumptions. 1242 A1 (Incoherence). $||A^*||_{2\to\infty} \leq \gamma \Delta^*$. 1243

A2 (Row- and column-wise independence). For any $m \in [n]$, the entries in the m^{th} row and column 1244 of A are independent from others, i.e. $\{A_{ij} : i = m \text{ or } j = m\} \perp \{A_{ij} : i \neq m, j \neq m\}$. 1245

A3 (Spectral norm concentration). For some $\delta_0 \in (0, 1)$, suppose $\mathbb{P}(||A - A^*||_2 \le \gamma \Delta^*) \ge 1 - \delta_0$. 1246 A4 (Row concentration). Suppose $\varphi(x)$ is continuous and non-decreasing in \mathbb{R}_+ and $\varphi(x)/x$ is 1247 non-increasing for x > 0. Additionally $\varphi(0) = 0$ and $32\kappa \max\{\gamma, \varphi(\gamma)\} \le 1$. Let there be some 1248 $\delta_1 \in (0,1)$ such that for any $m \in [n]$ and $w \in \mathbb{R}^n$ 1249

1250 1251

1252

1253

$$\mathbb{P}\left(|(A - A^*)_{m \cdot} w| \le \Delta^* \|w\|_{\infty} \varphi\left(\frac{\|w\|_2}{\sqrt{n} \|w\|_{\infty}}\right)\right) \ge 1 - \frac{\delta_1}{n}.$$

Lemma E.1 (Theorem 2.1 Abbe et al. (2020)). Under Assumptions 1 to 4, with probability at least $1-\delta_0-2\delta_1$, we have

$$\min_{s \in \{\pm 1\}} \left\| su_i - \frac{Au_i^*}{\lambda_i^*} \right\|_{\infty} \lesssim \kappa(\kappa + \varphi(1))(\gamma + \varphi(\gamma)) \left\| u_i^* \right\|_{\infty} + \frac{\gamma \left\| A^* \right\|_{2 \to \infty}}{\Delta^*}.$$

E.1 ENTRYWISE ANALYSIS FOR EIGENVECTORS FOR ROS.

1259 In this subsection, we will show that the top eigenvector of A sampled from ROS exhibits the entrywise behavior discussed above. More formally, we show the following lemma. 1261

Lemma E.2. Fix $\rho \in (0,1)$ and $a, b \in \mathbb{R}$ such that $\max\{|a|, |b|\} > 0$. Let $(A, \sigma^*) \sim \mathsf{ROS}_n(\rho, a, b)$. 1262 Condition on σ^* satisfying E from (9). Let $A^* := \mathbb{E}[A \mid \sigma^*]$. Define (λ_1, u_1) and (λ_1^*, u_1^*) as above. 1263 Then with probability 1 - o(1)1264

1278 1279 1280

$$\min_{s \in \{\pm 1\}} \left\| su_1 - \frac{Au_1^*}{\lambda_1^*} \right\|_{\infty} \le \frac{C}{\sqrt{n \log n}}$$

for some constant $C := C(\rho, a, b) > 0$. 1267

1268 According to the definition of the ROS model, we have $A = \operatorname{zd}\left(\sqrt{\frac{\log n}{n}} v^* v^{*\top} + W\right)$. The entire 1269 1270 analysis is done conditioned on σ^* , so the only randomness in this analysis is from the added noise 1271 matrix W. We verify Assumptions 1-4 required to apply Lemma E.1 using similar ideas as (Abbe 1272 et al., 2020, Theorem 3.1).

1273 First, observe that $A^* = \mathsf{zd}(v^* v^* \sqrt{\log n/n})$. Let (λ_1^*, u_1^*) be the top eigenpair. The corresponding 1274 eigengap quantity defined in (13) is $\Delta^* := |\lambda_1^*| \wedge \min_{2 \le i \le n} |\lambda_1^* - \lambda_i^*|$. We begin by characterizing $u_1^*, \lambda_1^*, \text{ and } \Delta^*.$ 1276

Lemma E.3. Let (λ_1^*, u_1^*) be the top eigenpair of A^* . Then 1277

$$u_1^* = \frac{(1+o(1))v^*}{\|v^*\|_2} \quad \lambda_1^* = (1+o(1))\sqrt{\frac{\log n}{n}} \|v^*\|_2^2, \quad \Delta^* \approx |\lambda_1^*| = \Theta(\sqrt{n\log n}).$$

Proof. Note that $v^* v^* \sqrt{\log n/n}$ is a rank-1 matrix. Let $|\tilde{\lambda}_1| \ge \cdots \ge |\tilde{\lambda}_n|$ be its eigenvalues. Then 1281 1282 we have that only non-zero eigenvalue is $\tilde{\lambda}_1 = \sqrt{\frac{\log n}{n}} \|v^*\|_2^2 = \Theta(\sqrt{n \log n})$ and the corresponding 1283 normalized eigenvector is $v^* / \|v^*\|_2$ and $\tilde{\lambda}_2 = \cdots = \tilde{\lambda}_n = 0$. After zeroing out the diagonal, the 1284 entries of the corresponding eigenvector $v^* / ||v^*||_2$ will be perturbed by a factor of (1 + o(1)) since 1285 the diagonal correction is of the order of $O(\sqrt{\log n/n})$. Hence, we obtain $u_1^* = (1+o(1))v^*/\|v^*\|_2$. 1286 By Weyl's inequality, we calculate the effect of zeroing out the diagonal on the eigenvalue: 1287

$$|\lambda_1^* - \tilde{\lambda}_1| \le \left\| v^* {v^*}^\top \sqrt{\log n/n} - \mathsf{zd}\left(v^* {v^*}^\top \sqrt{\log n/n} \right) \right\|_2 = O\left(\sqrt{\frac{\log n}{n}}\right).$$

Therefore, 1291

1290

1292
1293
$$\lambda_1^* = \tilde{\lambda}_1 + O(\sqrt{\log n/n}) = \sqrt{\frac{\log n}{n}} \|v^*\|_2^2 + O(\sqrt{\log n/n}) = (1 + o(1))\sqrt{\frac{\log n}{n}} \|v^*\|_2^2 \asymp \sqrt{n \log n}$$
1294

Applying Weyl's inequality, for $2 \le i \le n$, we get $|\lambda_i^*| = O(\sqrt{\log n/n})$. Hence, $\Delta^* \approx |\lambda_1^*| \asymp$ 1295 $\sqrt{n \log n}$.

Proof of Lemma E.2. We will let $\gamma := \frac{3\sqrt{n}}{\Delta^*} = 1/\Theta(\sqrt{\log n})$, due to Lemma E.3. Let us now verify Assumption 1. For any $i \in C_+$

$$\|A_{i\cdot}^*\|_2 = \sqrt{|C_+^{-i}| \cdot a^4 \frac{\log n}{n} + |C_-| \cdot a^2 b^2 \frac{\log n}{n}}$$

$$=\sqrt{(1+o(1))\rho n \cdot a^4 \frac{\log n}{n} + (1+o(1))(1-\rho)n \cdot a^2 b^2 \frac{\log n}{n}} = \Theta(\sqrt{\log n}),$$

where the second step follows from using Lemma B.3. Similarly, also for any $i \in C_{-}$

$$\|A_{i\cdot}^*\|_2 = \sqrt{|C_+| \cdot a^2 b^2 \frac{\log n}{n} + |C_-^{-i}| \cdot b^4 \frac{\log n}{n}} = \Theta(\sqrt{\log n}).$$

Overall, combining these two we obtain $||A^*||_{2\to\infty} = \Theta(\sqrt{\log n}) \le 3\sqrt{n} = \gamma \Delta^*$, verifying As-sumption 1. Assumption 2 on row and column-wise independence trivially holds due to the i.i.d. noise matrix W (up to symmetry).

To verify Assumption 3 on spectral norm concentration, first observe that $A - A^* = W$, where W is the zero diagonal symmetric matrix with i.i.d. $\mathcal{N}(0,1)$ entries. Applying (Bandeira et al., 2017, Proposition 3.3), we have that with probability at least $1 - e^{-n/2}$,

$$||A - A^*||_2 = ||W||_2 \le 3\sqrt{n} = \gamma \Delta^*$$

Therefore, Assumption 3 holds with $\delta_0 = e^{-n/2}$. We now turn our attention to Assumption 4. Let us choose $\varphi(x) = cx$ for some constant c > 0 which we will decide later. Clearly, φ is continuous, non-decreasing in \mathbb{R}_+ with $\varphi(0) = 0$, and $\varphi(x)/x = c$ is also non-increasing in $(0, \infty)$. Letting $\kappa = 1$, it is straightforward to see that $32\kappa \max\{\gamma, \varphi(\gamma)\} = 32\gamma \max\{1, c\} = o(1) \leq 1$, as $\gamma = o(1).$

We now verify the row concentration part of the assumption. Using Lemma E.3, we have $\Delta^* \approx$ $|\lambda_1^*| \geq \max\{\rho a^2, (1-\rho)b^2\}\sqrt{n\log n}$. Therefore, it holds that for any $\epsilon > 0$, there is a sufficiently large n such that $\Delta^* \ge (1-\epsilon) \max\{\rho a^2, (1-\rho)b^2\} \sqrt{n \log n}$. Moreover, for any fixed $w \in \mathbb{R}^n$, one can say that

$$\begin{aligned} \Delta^* \|w\|_{\infty} \varphi\left(\frac{\|w\|_2}{\sqrt{n} \|w\|_{\infty}}\right) &\geq \frac{(1-\epsilon) \max\{\rho a^2, (1-\rho)b^2\} \sqrt{n\log n} \cdot \|w\|_{\infty} \cdot c \|w\|_2}{\sqrt{n} \|w\|_{\infty}} \\ &= (1-\epsilon) c \max\{\rho a^2, (1-\rho)b^2\} \sqrt{\log n} \|w\|_2. \end{aligned}$$

Additionally, for any fixed $m \in [n]$, $(A - A^*)_m w \sim \mathcal{N}(0, ||w_{-m}||_2^2)$. Therefore,

$$\mathbb{P}\left(\left|(A - A^*)_{m \cdot} w\right| \le \Delta^* \left\|w\right\|_{\infty} \varphi\left(\frac{\|w\|_2}{\sqrt{n} \left\|w\right\|_{\infty}}\right)\right)$$

$$\geq \mathbb{P}\left(\left|(A - A^*)_{m} w\right| \le (1 - \epsilon)c \max\{\rho a^2, (1 - \rho)b^2\} \sqrt{\log n} \|w\|_2\right)$$

$$= \mathbb{P}\left(|\mathcal{N}(0, \|w_{-m}\|_{2}^{2})| \le (1-\epsilon)c \max\{\rho a^{2}, (1-\rho)b^{2}\} \sqrt{\log n} \|w\|_{2} \right)$$

$$\geq \mathbb{P}\left(|\mathcal{N}(0, \|w\|_{2}^{2})| \leq (1-\epsilon)c \max\{\rho a^{2}, (1-\rho)b^{2}\}\sqrt{\log n} \|w\|_{2}\right)$$

1339
1340
$$= \mathbb{P}\left(|\mathcal{N}(0,1)| \le (1-\epsilon)c \max\{\rho a^2, (1-\rho)b^2\}\sqrt{\log n}\right)$$

1340
$$= \mathbb{P}\left(|\mathcal{N}(0,1)| \le (1-\epsilon)c \max\{\rho a^2, (1-\rho)b^2\} \sqrt{\log n}\right)$$
1341

$$\geq 1 - \frac{2e^{-(1-\epsilon)^2 c^2 \max\{\rho a^2, (1-\rho)b^2\}^2 \log n/2}}{(1-\epsilon)c \max\{\rho a^2, (1-\rho)b^2\}\sqrt{\log n}\sqrt{2\pi}}$$

(using Lemma B.2)

$$\frac{(1-\epsilon)c\max\{\rho a^2, (1-\rho)b^2\}}{2n^{-(1-\epsilon)^2c^2}\max\{\rho a^2, (1-\rho)b^2\}}$$

$$= 1 - \frac{1}{(1-\epsilon)c \max\{\rho a^2, (1-\rho)b^2\}\sqrt{2\pi \log n}}$$

Therefore, letting

1348
1349
$$\delta_1 = \frac{2n^{1-(1-\epsilon)^2 c^2 \max\{\rho a^2, (1-\rho)b^2\}^2/2}}{(1-\epsilon)c \max\{\rho a^2, (1-\rho)b^2\}\sqrt{2\pi \log n}}, \text{ and setting } c = \frac{2}{(1-\epsilon)\max\{\rho a^2, (1-\rho)b^2\}},$$

we get $\delta_1 = o(1)$. Finally, applying Lemma E.1, we obtain that with probability $1 - \delta_0 - 2\delta_1 = 1 - o(1)$

1355 1356

$$\leq (1+c)(1+c)\gamma \left\|u_{1}^{*}\right\|_{\infty} + \frac{\gamma \left\|A^{*}\right\|_{2 \to \infty}}{\Delta^{*}} \quad (\text{since } \kappa = 1 \text{ and } \varphi(x) = cx)$$

1357 1358

$$= \frac{1}{\Theta(\sqrt{n\log n})}.$$

1359 1360 1361

1363

1364

1365

1367

1375 1376 1377

We used $\gamma = \frac{1}{\Theta(\sqrt{\log n})}$, $\|u_1^*\|_{\infty} = O(\frac{1}{\sqrt{n}})$, $\|A^*\|_{2 \to \infty} = \sqrt{\log n}$, and $\Delta^* = \sqrt{n \log n}$.

E.2 ENTRYWISE ANALYSIS OF EIGENVECTORS FOR SBM.

In this subsection, we show that the similar behavior also holds for eigenvectors of A sampled from the SBM. More specifically, we restrict ourselves to the case when the expectation A^* (after the appropriate diagonal correction) has rank 2. This is achieved when

$$\frac{a_1}{b} \neq \frac{b}{a_2}$$

In this case, the eigenvectors that correspond to the top two leading eigenvalues (in magnitude) exhibit the entrywise behavior, which is formalized in the following lemma.

1371 1372 1373 1374 Lemma E.4. Let $\rho \in (0,1)$ and $a_1, a_2, b > 0$ such that $a_1a_2 \neq b^2$. Let $(A, \sigma^*) \sim$ SBM_n(ρ, a_1, a_2, b). Condition on σ^* such that the event E from (9) holds and let $A := \mathbb{E}[A \mid \sigma^*]$. Define $\{(\lambda_i, u_i)\}_{i \in [n]}$ and $\{(\lambda_i^*, u_i^*)\}_{i \in [n]}$ as above. Then with probability $1 - O(n^{-3})$

$$\min_{s_1 \in \{\pm 1\}} \left\| s_1 u_1 - \frac{A u_1^*}{\lambda_1^*} \right\|_{\infty} \le \frac{C}{\sqrt{n} \log \log n} \quad and \quad \min_{s_2 \in \{\pm 1\}} \left\| s_2 u_2 - \frac{A u_2^*}{\lambda_2^*} \right\|_{\infty} \le \frac{C}{\sqrt{n} \log \log n},$$

1378 for some constant $C := C(\rho, a_1, a_2, b)$.

This again requires verifying Assumptions 1-4 for the top two eigenpairs. We note that Dhara et al.
 (2022b) showed a similar lemma for the special case of Planted Dense Subgraph (PDS), and our proof just generalizes their results. In order to do this, we first note down a couple of important lemmas. The first one directly establishes the spectral norm concentration (Assumption 3).

Lemma E.5. Let $\rho \in (0, 1)$ and $a_1, a_2, b > 0$. Sample $(A, \sigma^*) \sim \text{SBM}_n(\rho, a_1, a_2, b)$. Condition on σ^* such that the event E from (9) holds. Let $A^* := \mathbb{E}[A \mid \sigma^*]$, then there exists a constant $c_1 = c_1(\rho, a_1, a_2, b) > 0$ such that

 $\mathbb{P}(\|A - A^*\|_2 \le c_1 \sqrt{\log n}) \ge 1 - n^{-3}.$

1388 1389 *Proof.* The lemma is a special case of (Hajek et al., 2016, Theorem 5), invoking the theorem with c = 3.

The next lemma establishes that the leading two eigenvalues of A^* , in the rank-2 case, are different in the following sense.

Lemma E.6. Consider $\rho \in (0,1)$ and $a_1, a_2, b > 0$ such that $a_1a_2 \neq b^2$. Let $(A, \sigma^*) \sim SBM_n(\rho, a_1, a_2, b)$. Condition on a labelling σ^* such that the event E from (9) holds. Let $A^* := \mathbb{E}[A \mid \sigma^*]$. Then the top two eigenvalues in magnitude are given by $\lambda_1^* = (1 + o(1))\theta_1 \log n$ and $\lambda_2^* = (1 + o(1))\theta_2 \log n$, for some non-zero constants $\theta_1 \neq \theta_2$ in terms of (ρ, a_1, a_2, b) . As a consequence, $|\lambda_1^* - \lambda_2^*| = \Theta(\log n)$.

1399 *Proof.* The proof follows similar arguments as the proof of (Dhara et al., 2022b, Lemma 3.2). We 1400 note that they prove the special case of the PDS when $a_2 = b$, but the same argument directly 1401 generalizes as long as $a_1a_2 \neq b^2$.

1402

1403 Proof of Lemma E.4. The entire analysis is done conditioned on σ^* such that E holds. We will verify Assumptions 1-4 for the leading two eigenpairs. First note that after adding a diagonal matrix

D, whose entries are $O(\frac{\log n}{n})$, the matrix $A^* + D$ has rank 2, and its remaining eigenvalues satisfy $\lambda_3 = \cdots = \lambda_n = 0$, where $|\lambda_1| \ge \ldots |\lambda_n|$. Applying Weyl's inequality for $3 \le i \le n$,

$$\lambda_i^* - \tilde{\lambda}_i | = |\lambda_i^*| \le ||D||_2 = O(\log n/n).$$

$$\tag{14}$$

By the definition of the eigengap quantity in (13) for both the eigenvalues respectively

$$\Delta_1^* := |\lambda_1^*| \wedge \min_{i \neq 1} |\lambda_i^* - \lambda_1^*| = \Theta(\log n) \text{ and } \Delta_2^* := |\lambda_2^*| \wedge \min_{i \neq 2} |\lambda_i^* - \lambda_2^*| = \Theta(\log n),$$

where we used (14) and Lemma E.6. We also define $\kappa_1 := \frac{|\lambda_1^*|}{\Delta_1^*}$ and $\kappa_2 := \frac{|\lambda_2^*|}{\Delta_2^*}$. We first make an inportant observation, that to verify Assumptions 1-4 for both eigenpairs separately, it suffices to just verify them with Δ^* and κ such that

$$\Delta^* := \min\{\Delta_1^*, \Delta_2^*\} = \Theta(\log n) \text{ and } \kappa := \max\{\kappa_1, \kappa_2\}$$

To verify Assumption 1, first let $\tau = 2 \max\{a_1, a_2, b\}$. Fixing any $i \in C_+$,

$$||A_{i\cdot}||_{2} = \sqrt{|C_{+}^{-i}|\left(\frac{a_{1}\log n}{n}\right)^{2} + |C_{-}|\left(\frac{b\log n}{n}\right)^{2}}$$
$$= \sqrt{(1+o(1))n\left(\frac{a_{1}^{2}\log^{2}n}{n}\right) + (1+o(1))(1-n)\left(\frac{b^{2}\log^{2}n}{n}\right)} < \frac{\tau 1}{2}$$

$$= \sqrt{(1+o(1))\rho\left(\frac{a_1^2\log^2 n}{n}\right) + (1+o(1))(1-\rho)\left(\frac{b^2\log^2 n}{n}\right)} \le \frac{\tau\log n}{\sqrt{n}}.$$

Similarly, even for $i \in C_{-}$

$$\|A_{i\cdot}\|_{2} = \sqrt{|C_{+}|\left(\frac{b\log n}{n}\right)^{2} + |C_{-}^{-i}|\left(\frac{a_{2}\log n}{n}\right)^{2}} \le \frac{\tau\log n}{\sqrt{n}}$$

Combining both bounds, we obtain

$$\|A^*\|_{2\to\infty} \le \frac{\tau \log n}{\sqrt{n}}.$$

We now define the parameter γ in terms of Δ^* and the constant c_1 from Lemma E.5:

$$\gamma \triangleq \frac{c_1 \sqrt{\log n}}{\Delta^*} = \frac{c_1 \sqrt{\log n}}{\Theta(\log n)} = o(1)$$

Then $\gamma \Delta^* = c_1 \sqrt{\log n} = \Omega(\sqrt{\log n})$, which dominates $\tau \log n / \sqrt{n}$. This implies $||A^*||_{2 \to \infty} \leq \gamma \Delta^*$, verifying Assumption 1. Assumption 2 trivially holds due to the conditional independence of the entries of A, conditioned on σ^* . By Lemma E.5, Assumption 3 holds with $\delta_0 = n^{-3}$

$$\mathbb{P}(\|A - A^*\|_2 \le \gamma \Delta^*) \ge 1 - n^{-3}$$

To verify Assumption 4, we let

$$\varphi(x) \triangleq \frac{(2\tau+4)\log n}{\Delta^*(1 \vee \log(1/x))} \text{ for } x > 0 \text{ and } \varphi(0) = 0.$$

It is straightforward to verify that φ satisfies the desired property stated in Assumption 4 and $\varphi(\gamma) =$ $O(1/\log \log n)$. Also, $\kappa = O(1)$ since both $\Delta_1^* \simeq \Delta_2^* \simeq \log n$, and by Lemma E.6, also $|\lambda_1^*| \simeq |\lambda_2^*| \simeq 0$ $\log n$. This implies $32\kappa \max\{\gamma, \varphi(\gamma)\} = o(1)$ verifying the first part of the assumption.

To verify the row concentration part, we simply apply (Abbe et al., 2020, Lemma 7) with p = $\tau \log n/n$ and $\alpha = 4/\tau$. We obtain that for a fixed vector $w \in \mathbb{R}^n$ and $m \in [n]$,

1451
1452
1453
1454
$$\mathbb{P}\left(|(A-A^*)_m \cdot w| \le \frac{(2\tau+4)\log n}{\max\left\{1, \log\left(\frac{\sqrt{n}\|w\|_{\infty}}{\|w\|_2}\right)\right\}} \|w\|_{\infty}\right) \ge 1 - 2n^{-4}.$$

Substituting the definition of Δ^* and $\varphi(\cdot)$,

1457
$$\mathbb{P}\left(\left|(A-A^*)_m \cdot w\right| \le \Delta^* \left\|w\right\|_{\infty} \varphi\left(\frac{\|w\|_2}{\sqrt{n} \left\|w\right\|_{\infty}}\right)\right) \ge 1 - 2n^{-4}$$

which verifies Assumption 4 with $\delta_1 = 2n^{-3}$. Finally, applying Lemma E.1, with probability $1 - \delta_0 - 2\delta_1 = 1 - O(n^{-3}),$

1465
1466 We used
$$\gamma = \frac{1}{\Theta(\sqrt{\log n})}, \varphi(\gamma) = O(\frac{1}{\log \log n}), \|u_1^*\|_{\infty} = O(\frac{1}{\sqrt{n}}), \|A^*\|_{2 \to \infty} = O\left(\frac{\log n}{\sqrt{n}}\right), \text{ and } \Delta^* = \Theta(\log n).$$

Similarly, with probability $1 - O(n^{-3})$, we have

$$\min_{s_2 \in \{\pm 1\}} \left\| s_2 u_2 - \frac{A u_2^*}{\lambda_2^*} \right\|_{\infty} = O\left(\frac{1}{\sqrt{n} \log \log n}\right)$$

The proof is complete by a union bound.

F PROOFS AND ALGORITHMS FOR ROS.

/

In Appendix F.1, we first derive the form of genie scores. In Appendices F.2 we give our spectral algorithm formally and prove Theorem 1. Finally, in Appendix F.3, we provide degree-profiling algorithm under enormous BEC and BSC channel.

F.1 GENIE SCORES' FORMULA WHEN NO SIDE INFORMATION

We start by noting the form of genie scores when no side information is present.

Lemma F.1. Fix $\rho \in (0,1)$ and $a, b \in \mathbb{R}$ such that $\max\{|a|, |b|\} > 0$. Let $(A, \sigma^*) \sim \mathsf{ROS}_n(\rho, a, b)$. Then for any $i \in [n]$

、

$$\begin{array}{ll} & 1492 \\ 1493 \\ 1494 \\ 1494 \\ 1495 \\ 1496 \\ 1496 \\ 1497 \end{array} \\ & x_i^* = (a-b)\sqrt{\frac{\log n}{n}} \left(a \sum_{j \in C_+^{-i}} A_{ij} + b \sum_{j \in C_-^{-i}} A_{ij} \right) + \frac{\log n}{2n} (|C_+^{-i}| (a^2b^2 - a^4) + |C_-^{-i}| (b^4 - a^2b^2)) \\ + \log \left(\frac{\rho}{1-\rho} \right). \end{array}$$

 Moreover, conditioned on the event E from (9), the genie score vector $z^* \in \mathbb{R}^n$ can be written as

$$z^* = (a-b)\sqrt{\frac{\log n}{n}}Av^* + \left(\gamma + \log\left(\frac{\rho}{1-\rho}\right)\right)\mathbf{1}_n + o(1),$$

where $\gamma = (\rho(a^2b^2 - a^4) + (1 - \rho)(b^4 - a^2b^2)) \log n/2$ and v^* is given by (2).

Proof. For ease of notation, we denote $f(n) = \sqrt{\log n/n}$. First, note that ROS is a special case of GBM with $\mathcal{P}_+ = \mathcal{N}(a^2 f(n), 1), \mathcal{P}_- = \mathcal{N}(b^2 f(n), 1)$ and $\mathcal{Q} = \mathcal{N}(abf(n), 1)$. Applying Lemma

4.2 for this special case, we obtain the Genie score expressions; for any $i \in [n]$, $z_i^* = \sum_{j \in C_i^{-i}} \log \left(e^{-\frac{(A_{ij} - a^2 f(n))^2}{2}} \middle/ e^{-\frac{(A_{ij} - abf(n))^2}{2}} \right)$ $+ \sum_{i \in C^{-i}} \log \left(e^{-\frac{(A_{ij} - abf(n))^2}{2}} \middle/ e^{-\frac{(A_{ij} - b^2 f(n))^2}{2}} \right) + \log \left(\frac{\rho}{1 - \rho} \right)$ $=\sum_{j\in C^{-i}}\left(\frac{(A_{ij}-abf(n))^2}{2}-\frac{(A_{ij}-a^2f(n))^2}{2}\right)$ + $\sum_{i=1,\dots,n} \left(\frac{(A_{ij} - b^2 f(n))^2}{2} - \frac{(A_{ij} - abf(n))^2}{2} \right) + \log\left(\frac{\rho}{1-\rho}\right)$ $= (a^{2} - ab)f(n) \sum_{i \in C^{-i}} A_{ij} + (ab - b^{2})f(n) \sum_{i \in C^{-i}} A_{ij}$ $+\frac{f^{2}(n)}{2}(|C_{+}^{-i}|(a^{2}b^{2}-a^{4})+|C_{-}^{-i}|(b^{4}-a^{2}b^{2}))+\log\left(\frac{\rho}{1-\rho}\right)$ $= (a-b)\sqrt{\frac{\log n}{n}} \left(a \sum_{i \in C_{+}^{-i}} A_{ij} + b \sum_{j \in C_{-}^{-i}} A_{ij} \right)$ $+\frac{\log n}{2n}(|C_{+}^{-i}|(a^{2}b^{2}-a^{4})+|C_{-}^{-i}|(b^{4}-a^{2}b^{2}))+\log\left(\frac{\rho}{1-\rho}\right).$ (15)Conditioned on E, simplifying a term from (15)

$$\begin{aligned} \frac{\log n}{2n} (|C_{+}^{-i}|(a^{2}b^{2} - a^{4}) + |C_{-}^{-i}|(b^{4} - a^{2}b^{2})) \\ &= \frac{\log n}{2n} \left(\rho n(a^{2}b^{2} - a^{4}) + (1 - \rho)n \cdot (b^{4} - a^{2}b^{2}) \right) + o(1) \\ &= \left(\rho (a^{2}b^{2} - a^{4}) + (1 - \rho)(b^{4} - a^{2}b^{2}) \right) \log n/2 + o(1) \\ &= \gamma + o(1). \end{aligned}$$

Therefore, substituting this in (15), we obtain that for any $i \in [n]$:

$$z_{i}^{*} = (a-b)\sqrt{\frac{\log n}{n}} \left(a \sum_{j \in C_{+}^{-i}} A_{ij} + b \sum_{j \in C_{-}^{-i}} A_{ij} \right) + \gamma + \log\left(\frac{\rho}{1-\rho}\right) + o(1).$$

$$= (a-b)\sqrt{\frac{\log n}{n}}A_i \cdot v^* + \gamma + \log\left(\frac{\rho}{1-\rho}\right) + o(1).$$

1553 Writing the same for all $i \in [n]$ in vector notation, we obtain

$$z^* = (a-b)\sqrt{\frac{\log n}{n}}Av^* + \left(\gamma + \log\left(\frac{\rho}{1-\rho}\right)\right)\mathbf{1}_n + o(1). \quad \Box$$

Roughly speaking, these genie scores in absolute value are on a $\log n$ scale for both ROS and SBM. This is when the exact recovery becomes statistically possible and explains the scaling choices for both models.

F.2 SPECTRAL ALGORITHM AND PROOF OF THEOREM 1

Below is our spectral algorithm which takes A (and optionally y when available) as input along with the parameters and returns an estimator $\hat{\sigma}_{spec}$. One of the (two) score vectors formed by the algorithm approximates the genie score z^* .

 Imput: An <i>n</i>×<i>n</i> observation matrix <i>A</i> and parameters (<i>ρ</i>, <i>a</i>, <i>b</i>). Optionally, side information <i>y</i> ∈ such that we can compute likelihoods of laws <i>S</i>₊ and <i>S</i> Output: An estimate of community assignments ô_{spec}. 1: Compute leading eigenpairs. Compute the top eigenpair of <i>A</i>, denoted by (<i>λ</i>₁, <i>u</i>₁), whth <i>λ</i>₁ ≥ ··· ≥ <i>λ</i>_n . 2: Compute coefficients of linear combination. <i>c</i>₁ := √<i>n</i> log <i>n</i>·(<i>a</i>-<i>b</i>)·(<i>ρa</i>²+(1-<i>ρ</i>)<i>b</i>²)^{3/2} and <i>γ</i> := (<i>ρ</i>(<i>a</i>²<i>b</i>² - <i>a</i>⁴) + (1 - <i>ρ</i>)(<i>b</i>⁴ - <i>a</i>²<i>b</i>²)) 3: Compute spectral scores. For any <i>s</i> ∈ {±1}, prepare the spectral score vectors as follows. No side information: <i>z</i>^(s) = <i>sc</i>₁<i>u</i>₁ + <i>γ</i>1_{<i>n</i>}. • Side information: <i>z</i>^(s) = <i>sc</i>₁<i>u</i>₁ + <i>γ</i>1_{<i>n</i>} + log (<i>S</i>₊(<i>y</i>)) 4: Remove sign ambiguity. For each <i>s</i> ∈ {±1}, let ô^(s) = sgn(<i>z</i>^(s)). • No side information: Return ô_{spec} = arg max_{{ô^(c)}:<i>s</i>∈{±1}} P(<i>σ</i>[*] = ô^(s) <i>A</i>). • BEC or BSC side information: Return ô_{spec} = arg max_{{ô^(c)}:<i>s</i>∈{±1}} P(<i>σ</i>[*] = ô^(s) <i>A</i>. The values <i>c</i>₁ and <i>γ</i> are carefully designed to emulate the genie score. Since the eigenvectors only recovered up to a global direction flip, we need to keep both candidates in the algorithm. <i>G</i> of them is approximating the genie score well. Finally, whichever one has the higher poste probability is picked in step 4. To show the proof of the score approximation guarantee, we need following lemma, whose proof is included at the end of this subsection. Lemma F2. Fix <i>ρ</i> ∈ (0, 1) and <i>a</i>, <i>b</i> ∈ ℝ<i>such that</i> max{[<i>a</i>], <i>b</i>] > 0. Let (<i>A</i>, σ[*]) ~ ROS_n(<i>ρ</i>, <i>a Condition on σ[*] satisfying E from (9) holds for it. Then there exists a constant <i>c</i> := <i>c</i>(<i>ρ</i>, <i>a</i>, <i>b</i>) <i>s</i> that with probability 1 - <i>O</i>(<i>n</i>⁻³), the following event holds</i> <i>E</i>₁ := {√(^{log n}/_n <i>Av</i>* _∞ ≤ <i>c</i> log <i>n</i>}.
Output: An estimate of community assignments $\hat{\sigma}_{spec}$. 1: Compute leading eigenpairs. Compute the top eigenpair of A , denoted by (λ_1, u_1) , where $ \lambda_1 \geq \cdots \geq \lambda_n $. 2: Compute coefficients of linear combination. $c_1 := \sqrt{n} \log n \cdot (a-b) \cdot (\rho a^2 + (1-\rho)b^2)^{3/2}$ and $\gamma := (\rho(a^2b^2 - a^4) + (1-\rho)(b^4 - a^2b^2))$ 3: Compute spectral scores. For any $s \in \{\pm 1\}$, prepare the spectral score vectors as follows. • No side information: $z^{(s)} = sc_1u_1 + \gamma 1_n$. • Side information: $z^{(s)} = sc_1u_1 + \gamma 1_n + \log\left(\frac{S_+}{S}(y)\right)$ 4: Remove sign ambiguity. For each $s \in \{\pm 1\}$, let $\hat{\sigma}^{(s)} = \operatorname{sgn}(z^{(s)})$. • No side information: Return $\hat{\sigma}_{spec} = \arg \max_{\{\hat{\sigma}^{(s)}:s\in\{\pm 1\}\}} \mathbb{P}(\sigma^* = \hat{\sigma}^{(s)} \mid A)$. • BEC or BSC side information: Return $\hat{\sigma}_{spec} = \arg \max_{\{\hat{\sigma}^{(s)}:s\in\{\pm 1\}\}} \mathbb{P}(\sigma^* = \hat{\sigma}^{(s)} \mid A)$. The values c_1 and γ are carefully designed to emulate the gene is score. Since the eigenvectors only recovered up to a global direction flip, we need to keep both candidates in the algorithm. of them is approximating the genie score well. Finally, whichever one has the higher poste probability is picked in step 4. To show the proof of the score approximation guarantee, we need following lemma, whose proof is included at the end of this subsection. Lemma F.2. Fix $\rho \in (0, 1)$ and $a, b \in \mathbb{R}$ such that $\max\{ a , b \} > 0$. Let $(A, \sigma^*) \sim \operatorname{ROS}_n(\rho, a Condition on \sigma^*$ satisfying E from (9) holds for it. Then three exists a constant $c := c(\rho, a, b)$ so that with probability $1 - O(n^{-3})$, the following event holds $E_1 := \left\{ \sqrt{\frac{\log n}{n}} Av^* _{\infty} \le c \log n \right\}$. Below is our primary lemma which shows the spectral and genie score vector approximation in norm. Lemma F.3. Fix $\rho \in (0, 1)$ and a, b such that $\max\{ a , b \} > 0$. Let $(A, \sigma^*) \sim \operatorname{ROS}_n(\rho, a a)$ and condition on σ^* such that E from (9) holds. Optionally, let $y \sim \operatorname{SI}(\sigma^*, S_+, S)$ for the charal laws (S_+, S) . Let z^* and $z^{(s)}$ be the genie score and the spectral sco
1: Compute leading eigenpairs. Compute the top eigenpair of A , denoted by (λ_1, u_1) , where $ \lambda_1 \ge \dots \ge \lambda_n $. 2: Compute coefficients of linear combination. $c_1 := \sqrt{n} \log n \cdot (a-b) \cdot (\rho a^2 + (1-\rho)b^2)^{3/2}$ and $\gamma := (\rho(a^2b^2 - a^4) + (1-\rho)(b^4 - a^2b^2))$ 3: Compute spectral scores. For any $s \in \{\pm 1\}$, prepare the spectral score vectors as follows. • No side information: $z^{(s)} = sc_1u_1 + \gamma 1_n$. • Side information: $z^{(s)} = sc_1u_1 + \gamma 1_n + \log\left(\frac{S_+}{S}(y)\right)$ 4: Remove sign ambiguity. For each $s \in \{\pm 1\}$, let $\hat{\sigma}^{(s)} = \operatorname{sgn}(s^{(s)})$. • No side information: Return $\hat{\sigma}_{spec} = \arg \max_{\{\hat{\sigma}^{(s)}:s\in\{\pm 1\}\}} \mathbb{P}(\sigma^* = \hat{\sigma}^{(s)} \mid A)$. • BEC or BSC side information: Return $\hat{\sigma}_{spec} = \arg \max_{\{\hat{\sigma}^{(s)}:s\in\{\pm 1\}\}} \mathbb{P}(\sigma^* = \hat{\sigma}^{(s)} \mid A)$. The values c_1 and γ are carefully designed to emulate the gene score. Since the eigenvectors only recovered up to a global direction flip, we need to keep both candidates in the algorithm. of them is approximating the genie score well. Finally, whichever one has the higher poste probability is picked in step 4. To show the proof of the score approximation guarantee, we need following lemma, whose proof is included at the end of this subsection. Lemma F.2. Fix $\rho \in (0, 1)$ and $a, b \in \mathbb{R}$ such that $\max\{ a , b \} > 0$. Let $(A, \sigma^*) \sim \operatorname{ROS}_n(\rho, a Condition on \sigma^* satisfying E from (9) holds for it. Then there exists a constant c := c(\rho, a, b) is that with probability 1 - O(n^{-3}), the following event holdsE_1 := \left\{ \sqrt{\frac{\log n}{n}} \ Av^*\ _{\infty} \le c \log n \right\}.Below is our primary lemma which shows the spectral and genie score vector approximation in norm.Lemma F.3. Fix \rho \in (0, 1) and a, b such that \max\{ a , b \} > 0. Let (A, \sigma^*) \sim \operatorname{ROS}_n(\rho, a) and condition on \sigma^* such that E from (9) holds. Optionally, let y \sim \operatorname{Sl}(\sigma^*, S_+, S) for the chara laws (S_+, S). Let z^* and z^{(s)} be the genie score and the spectral score vectors respectively for coresponding model. Then with probability 1 - o(1)$
$c_{1} := \sqrt{n} \log n \cdot (a-b) \cdot (\rho a^{2} + (1-\rho)b^{2})^{3/2} \text{ and } \gamma := (\rho(a^{2}b^{2} - a^{4}) + (1-\rho)(b^{4} - a^{2}b^{2}))$ 3: Compute spectral scores. For any $s \in \{\pm 1\}$, prepare the spectral score vectors as follows. • No side information: $z^{(s)} = sc_{1}u_{1} + \gamma 1_{n}$. • Side information: $z^{(s)} = sc_{1}u_{1} + \gamma 1_{n} + \log\left(\frac{S_{+}}{S_{-}}(y)\right)$ 4: Remove sign ambiguity. For each $s \in \{\pm 1\}$, let $\hat{\sigma}^{(s)} = \operatorname{sgn}(z^{(s)})$. • No side information: Return $\hat{\sigma}_{\operatorname{spec}} = \arg \max_{\{\hat{\sigma}^{(s)}:s\in\{\pm 1\}\}} \mathbb{P}(\sigma^{*} = \hat{\sigma}^{(s)} \mid A)$. • BEC or BSC side information: Return $\hat{\sigma}_{\operatorname{spec}} = \arg \max_{\{\hat{\sigma}^{(s)}:s\in\{\pm 1\}\}} \mathbb{P}(\sigma^{*} = \hat{\sigma}^{(s)} \mid A)$. The values c_{1} and γ are carefully designed to emulate the genie score. Since the eigenvectors only recovered up to a global direction flip, we need to keep both candidates in the algorithm. Construction of them is approximating the genie score well. Finally, whichever one has the higher posted following lemma, whose proof is included at the end of this subsection. Lemma F.2. Fix $\rho \in (0, 1)$ and $a, b \in \mathbb{R}$ such that $\max\{ a , b \} > 0$. Let $(A, \sigma^{*}) \sim \operatorname{ROS}_{n}(\rho, a Condition on \sigma^{*} satisfying E from (9) holds for it. Then there exists a constant c := c(\rho, a, b) is that with probability 1 - O(n^{-3}), the following event holdsE_{1} := \left\{ \sqrt{\frac{\log n}{n}} \ Av^{*}\ _{\infty} \le c \log n \right\}.Below is our primary lemma which shows the spectral and genie score vector approximation in norm.Lemma F.3. Fix \rho \in (0, 1) and a, b such that \max\{ a , b \} > 0. Let (A, \sigma^{*}) \sim \operatorname{ROS}_{n}(\rho, a) and condition on \sigma^{*} such that E from (9) holds. Optionally, let y \sim \operatorname{Sl}(\sigma^{*}, S_{+}, S_{-}) for the charles and condition on \sigma^{*} such that E from (9) holds. Optionally, let y \sim \operatorname{Sl}(\sigma^{*}, S_{+}, S_{-}) for the charles and condition on \sigma^{*} such that E from (9) holds. Optionally, let y \sim \operatorname{Sl}(\sigma^{*}, S_{+}, S_{-}) for the charles and condition on \sigma^{*} such that E from (9) holds. Optionally, let y \sim \operatorname{Sl}(\sigma^{$
 3: Compute spectral scores. For any s ∈ {±1}, prepare the spectral score vectors as follows. No side information: z^(s) = sc₁u₁ + γ1_n. Side information: z^(s) = sc₁u₁ + γ1_n + log (S+/S(y)) 4: Remove sign ambiguity. For each s ∈ {±1}, let ô^(s) = sgn(z^(s)). No side information: Return ô_{spec} = arg max_{{∂^(s):s∈{±1}}} P(σ* = ô^(s) A). BEC or BSC side information: Return ô_{spec} = arg max_{{∂^(s):s∈{±1}}} P(σ* = ô^(s) A). The values c₁ and γ are carefully designed to emulate the genie score. Since the eigenvectors only recovered up to a global direction flip, we need to keep both candidates in the algorithm. Got them is approximating the genie score well. Finally, whichever one has the higher poste probability is picked in step 4. To show the proof of the score approximation guarantee, we need following lemma, whose proof is included at the end of this subsection. Lemma F.2. Fix ρ ∈ (0, 1) and a, b ∈ ℝ such that max{ a , b } > 0. Let (A, σ*) ~ ROS_n(ρ, a Condition on σ* satisfying E from (9) holds for it. Then there exists a constant c := c(ρ, a, b) s that with probability 1 − O(n⁻³), the following event holds E₁ := {√(log n)/n Av* _∞ ≤ c log n}. Below is our primary lemma which shows the spectral and genie score vector approximation in norm. Lemma F.3. Fix ρ ∈ (0, 1) and a, b such that max{ a , b } > 0. Let (A, σ*) ~ ROS_n(ρ, a Condition on σ* such that E from (9) holds. Optionally, let y ~ Sl(σ*, S₊, S₋) for the charl may [x*, S). Let z* and z^(s) be the genie score and the spectral score vectors respectively for corresponding model. Then with probability 1 − o(1).
 No side information: z^(s) = sc₁u₁ + γ1_n. Side information: z^(s) = sc₁u₁ + γ1_n + log (S₊/S₋(y)) 4: Remove sign ambiguity. For each s ∈ {±1}, let ô^(s) = sgn(z^(s)). No side information: Return ô_{spec} = arg max_{{ô^(s):s∈{±1}}} P(σ* = ô^(s) A). BEC or BSC side information: Return ô_{spec} = arg max_{{ô^(s):s∈{±1}}} P(σ* = ô^(s) A). The values c₁ and γ are carefully designed to emulate the genie score. Since the eigenvectors only recovered up to a global direction flip, we need to keep both candidates in the algorithm. Of them is approximating the genie score well. Finally, whichever one has the higher poste probability is picked in step 4. To show the proof of the score approximation guarantee, we need following lemma, whose proof is included at the end of this subsection. Lemma F2. Fix ρ ∈ (0, 1) and a, b ∈ ℝ such that max{ a , b } > 0. Let (A, σ*) ~ ROS_n(ρ, a Condition on σ* satisfying E from (9) holds for it. Then there exists a constant c := c(ρ, a, b) s that with probability 1 − O(n⁻³), the following event holds E₁ := { √(log n)/n Av[*] _∞ ≤ c log n }. Below is our primary lemma which shows the spectral and genie score vector approximation in norm. Lemma F3. Fix ρ ∈ (0, 1) and a, b such that max{ a , b } > 0. Let (A, σ*) ~ ROS_n(ρ, a and condition on σ* such that E from (9) holds. Optionally, let y ~ Sl(σ*, S+, S) for the charla laws (S+, S-). Let z* and z^(s) be the genie score and the spectral score vectors respectively for corresponding model. Then with probability 1 − o(1),
• Side information: $z^{(s)} = sc_1u_1 + \gamma 1_n + \log\left(\frac{S_+}{S}(y)\right)$ 4: Remove sign ambiguity. For each $s \in \{\pm 1\}$, let $\hat{\sigma}^{(s)} = \operatorname{sgn}(z^{(s)})$. • No side information: Return $\hat{\sigma}_{\operatorname{spec}} = \arg \max_{\{\hat{\sigma}^{(s)}:s \in \{\pm 1\}\}} \mathbb{P}(\sigma^* = \hat{\sigma}^{(s)} \mid A)$. • BEC or BSC side information: Return $\hat{\sigma}_{\operatorname{spec}} = \arg \max_{\{\hat{\sigma}^{(s)}:s \in \{\pm 1\}\}} \mathbb{P}(\sigma^* = \hat{\sigma}^{(s)} \mid A)$. The values c_1 and γ are carefully designed to emulate the genie score. Since the eigenvectors only recovered up to a global direction flip, we need to keep both candidates in the algorithm. Of them is approximating the genie score well. Finally, whichever one has the higher poste probability is picked in step 4. To show the proof of the score approximation guarantee, we need following lemma, whose proof is included at the end of this subsection. Lemma F2. Fix $\rho \in (0, 1)$ and $a, b \in \mathbb{R}$ such that $\max\{ a , b \} > 0$. Let $(A, \sigma^*) \sim \operatorname{ROS}_n(\rho, a Condition on \sigma^* satisfying E from (9) holds for it. Then there exists a constant c := c(\rho, a, b) s that with probability 1 - O(n^{-3}), the following event holdsE_1 := \left\{ \sqrt{\frac{\log n}{n}} \ Av^*\ _{\infty} \le c \log n \right\}.Below is our primary lemma which shows the spectral and genie score vector approximation in norm.Lemma F3. Fix \rho \in (0, 1) and a, b such that \max\{ a , b \} > 0. Let (A, \sigma^*) \sim \operatorname{ROS}_n(\rho, a a) and condition on \sigma^* such that E from (9) holds. Optionally, let y \sim \operatorname{SI}(\sigma^*, S_+, S) for the charlaws (S_+, S). Let z^* and z^{(s)} be the genie score and the spectral score vectors respectively for corresponding model. Then with probability 1 - o(1),\min \ z^* - z^{(s)}\ = o(\log n)$
$z^{(s)} = sc_1u_1 + \gamma 1_n + \log\left(\frac{S_+}{S}(y)\right)$ 4: Remove sign ambiguity. For each $s \in \{\pm 1\}$, let $\hat{\sigma}^{(s)} = \operatorname{sgn}(z^{(s)})$. • No side information: Return $\hat{\sigma}_{\operatorname{spec}} = \arg \max_{\{\hat{\sigma}^{(s)}: s \in \{\pm 1\}\}} \mathbb{P}(\sigma^* = \hat{\sigma}^{(s)} \mid A)$. • BEC or BSC side information: Return $\hat{\sigma}_{\operatorname{spec}} = \arg \max_{\{\hat{\sigma}^{(s)}: s \in \{\pm 1\}\}} \mathbb{P}(\sigma^* = \hat{\sigma}^{(s)} \mid A)$. The values c_1 and γ are carefully designed to emulate the genie score. Since the eigenvectors only recovered up to a global direction flip, we need to keep both candidates in the algorithm. Of them is approximating the genie score well. Finally, whichever one has the higher poste probability is picked in step 4. To show the proof of the score approximation guarantee, we need following lemma, whose proof is included at the end of this subsection. Lemma F.2. Fix $\rho \in (0, 1)$ and $a, b \in \mathbb{R}$ such that $\max\{ a , b \} > 0$. Let $(A, \sigma^*) \sim \operatorname{ROS}_n(\rho, a Condition on \sigma^* satisfying E from (9) holds for it. Then there exists a constant c := c(\rho, a, b) is that with probability 1 - O(n^{-3}), the following event holdsE_1 := \left\{ \sqrt{\frac{\log n}{n}} \ Av^*\ _{\infty} \le c \log n \right\}.Below is our primary lemma which shows the spectral and genie score vector approximation in norm.Lemma F.3. Fix \rho \in (0, 1) and a, b such that \max\{ a , b \} > 0. Let (A, \sigma^*) \sim \operatorname{ROS}_n(\rho, a a) and condition on \sigma^* such that E from (9) holds. Optionally, let y \sim \operatorname{Sl}(\sigma^*, S_+, S) for the charlaws (S_+, S). Let z^* and z^{(s)} be the genie score and the spectral score vectors respectively for corresponding model. Then with probability 1 - o(1),min \ z^* - z^{(s)}\ = o(\log n)$
 4: Remove sign ambiguity. For each s ∈ {±1}, let ô^(s) = sgn(z^(s)). No side information: Return ô_{spec} = arg max_{{∂^(s):s∈{±1}}} P(σ* = ô^(s) A). BEC or BSC side information: Return ô_{spec} = arg max_{{∂^(s):s∈{±1}}} P(σ* = ô^(s) A). The values c₁ and γ are carefully designed to emulate the genie score. Since the eigenvectors only recovered up to a global direction flip, we need to keep both candidates in the algorithm. Of them is approximating the genie score well. Finally, whichever one has the higher poste probability is picked in step 4. To show the proof of the score approximation guarantee, we need following lemma, whose proof is included at the end of this subsection. Lemma F.2. Fix ρ ∈ (0, 1) and a, b ∈ ℝ such that max{ a , b } > 0. Let (A, σ*) ~ ROS_n(ρ, a Condition on σ* satisfying E from (9) holds for it. Then there exists a constant c := c(ρ, a, b) s that with probability 1 − O(n⁻³), the following event holds E₁ := {√(log n)/n Av[*] _∞ ≤ c log n}. Below is our primary lemma which shows the spectral and genie score vector approximation in norm. Lemma F.3. Fix ρ ∈ (0, 1) and a, b such that max{ a , b } > 0. Let (A, σ*) ~ ROS_n(ρ, a and condition on σ* such that E from (9) holds. Optionally, let y ~ Sl(σ*, S₊, S₋) for the charlaws (S₊, S₋). Let z* and z^(s) be the genie score and the spectral score vectors respectively for corresponding model. Then with probability 1 − o(1),
 No side information: Return ô_{spec} = arg max_{{∂^(s):s∈{±1}}} P(σ* = ô^(s) A). BEC or BSC side information: Return ô_{spec} = arg max_{{∂^(s):s∈{±1}}} P(σ* = ô^(s) A, The values c₁ and γ are carefully designed to emulate the genie score. Since the eigenvectors only recovered up to a global direction flip, we need to keep both candidates in the algorithm. Of them is approximating the genie score well. Finally, whichever one has the higher poste probability is picked in step 4. To show the proof of the score approximation guarantee, we need following lemma, whose proof is included at the end of this subsection. Lemma F.2. Fix ρ ∈ (0, 1) and a, b ∈ R such that max{ a , b } > 0. Let (A, σ*) ~ ROS_n(ρ, a Condition on σ* satisfying E from (9) holds for it. Then there exists a constant c := c(ρ, a, b) s that with probability 1 - O(n⁻³), the following event holds E₁ := {√(log n)/n Av* _∞ ≤ c log n}. Below is our primary lemma which shows the spectral and genie score vector approximation in norm. Lemma F.3. Fix ρ ∈ (0, 1) and a, b such that max{ a , b } > 0. Let (A, σ*) ~ ROS_n(ρ, a and condition on σ* such that E from (9) holds. Optionally, let y ~ Sl(σ*, S₊, S₋) for the charlaws (S₊, S₋). Let z* and z^(s) be the genie score and the spectral score vectors respectively for corresponding model. Then with probability 1 - o(1),
• BEC or BSC side information: Return $\hat{\sigma}_{spec} = \arg \max_{\{\hat{\sigma}^{(s)}:s \in \{\pm 1\}\}} \mathbb{P}(\sigma^* = \hat{\sigma}^{(s)} \mid A)$. The values c_1 and γ are carefully designed to emulate the genie score. Since the eigenvectors only recovered up to a global direction flip, we need to keep both candidates in the algorithm. Of them is approximating the genie score well. Finally, whichever one has the higher poster probability is picked in step 4. To show the proof of the score approximation guarantee, we need following lemma, whose proof is included at the end of this subsection. Lemma F.2. Fix $\rho \in (0, 1)$ and $a, b \in \mathbb{R}$ such that $\max\{ a , b \} > 0$. Let $(A, \sigma^*) \sim \operatorname{ROS}_n(\rho, a Condition on \sigma^* satisfying E from (9) holds for it. Then there exists a constant c := c(\rho, a, b) s that with probability 1 - O(n^{-3}), the following event holdsE_1 := \left\{ \sqrt{\frac{\log n}{n}} \ Av^*\ _{\infty} \le c \log n \right\}.$
The values c_1 and γ are carefully designed to emulate the genie score. Since the eigenvectors only recovered up to a global direction flip, we need to keep both candidates in the algorithm. Of of them is approximating the genie score well. Finally, whichever one has the higher poste probability is picked in step 4. To show the proof of the score approximation guarantee, we need following lemma, whose proof is included at the end of this subsection. Lemma F.2. Fix $\rho \in (0, 1)$ and $a, b \in \mathbb{R}$ such that $\max\{ a , b \} > 0$. Let $(A, \sigma^*) \sim \text{ROS}_n(\rho, a$ Condition on σ^* satisfying E from (9) holds for it. Then there exists a constant $c := c(\rho, a, b)$ s that with probability $1 - O(n^{-3})$, the following event holds $E_1 := \left\{ \sqrt{\frac{\log n}{n}} \ Av^*\ _{\infty} \le c \log n \right\}$. Below is our primary lemma which shows the spectral and genie score vector approximation in norm. Lemma F.3. Fix $\rho \in (0, 1)$ and a, b such that $\max\{ a , b \} > 0$. Let $(A, \sigma^*) \sim \text{ROS}_n(\rho, a$ and condition on σ^* such that E from (9) holds. Optionally, let $y \sim \text{SI}(\sigma^*, S_+, S)$ for the char laws (S_+, S) . Let z^* and $z^{(s)}$ be the genie score and the spectral score vectors respectively for corresponding model. Then with probability $1 - o(1)$, $\min \ z^* - z^{(s)}\ = o(\log n)$
$E_{1} := \left\{ \sqrt{\frac{\log n}{n}} \ Av^{*}\ _{\infty} \leq c \log n \right\}.$ Below is our primary lemma which shows the spectral and genie score vector approximation in norm. Lemma F.3. Fix $\rho \in (0,1)$ and a, b such that $\max\{ a , b \} > 0$. Let $(A, \sigma^{*}) \sim \operatorname{ROS}_{n}(\rho, c)$ and condition on σ^{*} such that E from (9) holds. Optionally, let $y \sim \operatorname{Sl}(\sigma^{*}, \mathcal{S}_{+}, \mathcal{S}_{-})$ for the charlaws $(\mathcal{S}_{+}, \mathcal{S}_{-})$. Let z^{*} and $z^{(s)}$ be the genie score and the spectral score vectors respectively for corresponding model. Then with probability $1 - o(1)$, $\min \ z^{*} - z^{(s)}\ = o(\log n)$
Below is our primary lemma which shows the spectral and genie score vector approximation in norm. Lemma F.3. Fix $\rho \in (0,1)$ and a, b such that $\max\{ a , b \} > 0$. Let $(A, \sigma^*) \sim \text{ROS}_n(\rho, \sigma)$ and condition on σ^* such that E from (9) holds. Optionally, let $y \sim \text{SI}(\sigma^*, S_+, S)$ for the charlaws (S_+, S) . Let z^* and $z^{(s)}$ be the genie score and the spectral score vectors respectively for corresponding model. Then with probability $1 - o(1)$, $\min \left\ z^* - z^{(s)} \right\ _{c} = o(\log n)$
Lemma F.3. Fix $\rho \in (0,1)$ and a, b such that $\max\{ a , b \} > 0$. Let $(A, \sigma^*) \sim \text{ROS}_n(\rho, c)$ and condition on σ^* such that E from (9) holds. Optionally, let $y \sim \text{SI}(\sigma^*, S_+, S)$ for the char laws (S_+, S) . Let z^* and $z^{(s)}$ be the genie score and the spectral score vectors respectively for corresponding model. Then with probability $1 - o(1)$, $\min \ z^* - z^{(s)}\ = o(\log n)$
$\min z^* - z^{(s)} = o(\log n)$
$s \in \{\pm 1\} \parallel^{\infty} \approx \parallel_{\infty} O(\log N).$
<i>Proof.</i> First of all, note that conditioned on E , we have $\lambda_1^* = (1 + o(1))\sqrt{\frac{\log n}{n}} \ v^*\ _2^2$, and u^*
$(1+o(1))\frac{v^{-}}{\ v^{*}\ _{2}}$. Additionally, $\ v^{*}\ _{2} = \sqrt{ C_{+} a^{2} + C_{-} b^{2}} = (1+o(1))\sqrt{n}\sqrt{\rho a^{2} + (1-\rho)^{2}}$. Using these, one can simplify
$\frac{c_1 A u_1^*}{\lambda_1^*} = \frac{(1+o(1))c_1 A v^*}{\lambda_1^* \ v^*\ _2} \approx \frac{\sqrt{n}\log n(a-b)(\rho a^2 + (1-\rho)b^2)^{3/2} A v^*}{\sqrt{\log n/n} \ v^*\ _2^3} \approx \sqrt{\frac{\log n}{n}}(a-b)A v^* + \frac{1}{n} \sum_{n=1}^{\infty} \frac{1}{n} \sum$
where in the last step we substitute $ v^* _{2}$.
Now the high probability event in this lemma is such that (i) the behavior of eigenvectors as st_i

in Lemma E.4 and (ii) the event E_1 from (16) hold. By Lemma E.2 and F.2, they both happen with probability 1 - o(1). Additionally, let s be the sign for which the conclusion of Lemma E.2 holds. We now analyze the three models separately.

• No side information: For every $i \in [n]$ 1621 $|z^{(s)} - z_i^*| = |c_1 s(u_1)_i + \gamma - z_i^*|,$ 1622 $= |c_1 s(u_1)_i + \gamma - z_{i+1}|$ = $\left| \frac{c_1 (Au_1^*)_i}{\lambda_1^*} + \gamma - z_i^* \right| + O\left(\frac{c_1}{\sqrt{n \log n}}\right)$ (by Lemma E.4 and the triangle inequality) (recall Algorithm 2) 1623 1624 1625 1626 $= \left| (1+o(1))\sqrt{\frac{\log n}{n}}(a-b)(Av^*)_i + \gamma - z_i^* \right| + O(\sqrt{\log n})$ 1627 1628 1629 (using (17) and $c_1 \asymp \sqrt{n} \log n$) $= \left| (1+o(1))\sqrt{\frac{\log n}{n}}(a-b)(Av^{*})_{i} + \gamma - \sqrt{\frac{\log n}{n}}(a-b)A_{i} \cdot v^{*} - \gamma - O(1) \right|$ 1633 $+O(\sqrt{\log n})$ (putting z_i^* from Lemma F.1) $= o(1)\sqrt{\frac{\log n}{n}} |(Av^*)_i| + O(\sqrt{\log n}) = o(\log n), \qquad (\text{recall } E_1 \text{ from (16)})$ 1635 1637 • Side information: By Lemma 4.2 and Algorithm 2 (step 3), when side information is provided, both the genie score vector z^* and the spectral score vector $z^{(s)}$ are achieved by 1639 adding $\ln\left(\frac{S_+}{S_-}(y)\right)$ to their counterpart when no side information is provided. Therefore, 1640 the triangle inequality along with the analysis in the no side information case gives us that 1641 for every $i \in [n]$, we have $|z_i^{(s)} - z_i^*| = o(\log n)$. 1642 1643 In either case, one can equivalently write conclusions for all $i \in [n]$ together in vector notation 1644 1645 $\min_{s \in \{\pm 1\}} \left\| z^{(s)} - z^* \right\|_{\infty} = o(\log n). \quad \Box$ 1646 1647 1648 We are finally set to prove our first main result in Theorem 1. 1650 1651 *Proof of Theorem 1.* We note that step 4 of Algorithm 2 keeps two candidates $\{\hat{\sigma}^{(s)} : s \in \{\pm 1\}\}$ 1652 and chooses the one which has maximum posterior probability. Therefore, to show that $\hat{\sigma}_{\text{spec}}$ achieves exact recovery above the IT threshold, it suffices to show that one of the two candidates 1654 achieves exact recovery, and $\hat{\sigma}_{MAP}$ also succeeds above the IT threshold, which ensures that the 1655 algorithm selects the correct vector by maximizing the posterior probability. The MAP estimator 1656 achieves exact recovery whenever $I^* >$ is already shown in Proposition 3.1. It remains to show that one of $\{\hat{\sigma}^{(s)} : s \in \{\pm 1\}\}$ succeeds. To this end, recall Lemma F.3 that with probability 1 - o(1), 1657 1658 $\min_{s \in \{\pm 1\}} \left\| z^* - z^{(s)} \right\|_{\infty} = o(\log n).$ 1659 Moreover, whenever $I^* > 1$, by Lemma D.1 and union bound over $i \in [n]$, there exists $\delta > 0$ such 1661 that 1662 $\mathbb{P}\left(\min_{i\in[n]}\sigma_i^* z_i^* > \delta \log n\right) = 1 - o(1).$ 1663 1664 Taking a union bound over these two events, there exists $\varsigma > 0$ and $s^* \in \{\pm 1\}$ such that 1665 1666

$$\mathbb{P}\left(\min_{i\in[n]}\sigma_i^* z_i^{(s^*)} > \varsigma \log n\right) = 1 - o(1).$$

1669 Since $\hat{\sigma}^{(s^*)} = \operatorname{sgn}(z^{(s^*)})$ in step 4, we obtain $\hat{\sigma}^{(s^*)}$ achieves exact recovery. As a consequence, even $\hat{\sigma}_{\text{spec}}$ achieves exact recovery above the IT threshold. In other words,

1672
$$\lim_{n \to \infty} \mathbb{P}(\hat{\sigma}_{\text{spec}} \text{ succeeds}) = 1. \quad \Box$$

1668

1673

We finally return to the proof of the lemma already mentioned.

1676 1677 Proof of Lemma F.2. For each $i \in [n]$, first define $Y_i = \sqrt{\frac{\log n}{n}} (Av^*)_i$. We first note that Y_i is 1678 a Gaussian random variable as it is the sum of at most n independent Gaussian random variables. 1679 Therefore, we have $Y_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$ for certain μ_i and σ_i^2 , which we will calculate later. Applying 1680 Lemma B.2 for any Y_i (after normalizing) yields

$$\mathbb{P}\left(\frac{|Y_i - \mu_i|}{\sigma_i} \ge 4\sqrt{\log n}\right) \le e^{-8\log n} = n^{-8}.$$

1683 1684 1685 Rearrangement of the terms using the triangle inequality along with a union bound over all $i \in [n]$ gives us

$$\mathbb{P}\left(\forall i \in [n] : |Y_i| \le |\mu_i| + 4\sigma_i \sqrt{\log n}\right) \ge 1 - n^{-7}.$$

1688 Therefore it simply suffices to show that for every $i \in [n]$, the quantity $|\mu_i| + 4\sigma_i \sqrt{\log n} = O(\log n)$. To this end, we first observe that $(Av^*)_i$ is the sum of n-1 independent Gaussian random variables 1690 all with means whose absolute values are $O\left(\sqrt{\frac{\log n}{n}}\right)$. Thus,

1692 1693

1694 1695

1699 1700 1701

1703 1704

1705

1682

1687

$$|\mu_i| = \sqrt{\frac{\log n}{n}} \mathbb{E}[(Av^*)_i] = \sqrt{\frac{\log n}{n}}(n-1) \cdot O\left(\sqrt{\frac{\log n}{n}}\right) = O(\log n).$$

1696 Similarly, $(Av^*)_i$ is the sum of n-1 independent Gaussian random variables with variances O(1). 1697 This gives us

$$\sigma_i^2 = \operatorname{Var}\left[\sqrt{\frac{\log n}{n}}(Av^*)_i\right] = \frac{\log n \cdot \operatorname{Var}\left[(Av^*)_i\right]}{n} = \frac{\log n \cdot O(n)}{n} = O(\log n)$$

which also implies

$$4\sigma_i \sqrt{\log n} = O(\log n).$$

1706 F.3 DEGREE-PROFILING ALGORITHM FOR BEC AND BSC CHANNELS

The following is a simple degree-profiling algorithm that tries to mimic the genie naïvely and achieves exact recovery if side information is substantial to shift the thresholds of exact recovery.

1710 Algorithm 3 Degree-Profiling algorithm for ROS in the presence of BEC or BSC side information.

1711 Input: An $n \times n$ observation matrix A and parameters (ρ, a, b) . The BEC side information y with 1712 parameter ϵ or BSC side information y with parameter α .

Output: An estimate of community assignments $\hat{\sigma}_{dp}$.

1: Let
$$S_+ := \{i : y_i = +1\}, S_- := \{i : y_i = -1\}$$
, and

1715 1716

1717

1718 1719

1722

1723

1724

$$\gamma := \left(\rho(a^2b^2 - a^4) + (1 - \rho)(b^4 - a^2b^2)\right)\log n/2$$

Compute $z \in \mathbb{R}^n$ such that, for every $i \in [n]$

$$z_i = a(a-b)\sqrt{\frac{\log n}{n}} \sum_{j \in S_+} A_{ij} + b(a-b)\sqrt{\frac{\log n}{n}} \sum_{j \in S_-} A_{ij} + \gamma.$$

2: Prepare the degree-profile score vector z^{dp} as follows.

• BEC side information: For any $i \in [n]$,

1725 1726 1727 $z_i^{dp} = \begin{cases} z_i & \text{if } y_i = 0; \\ +\infty, & \text{if } y_i = +1; \\ -\infty & \text{if } y_i = -1; \end{cases}$ • BSC side information:

$$z^{\rm dp} = z + \ln\left(\frac{1-\alpha}{\alpha}\right) y$$

1731 1732 3: Return $\hat{\sigma}_{dp} = \text{sgn}(z^{dp}).$

1728

1729 1730

1737

1738

1748 1749 1750

1753

1754

1755 1756 1757

1758

1759 1760

1761

1762

1763

1764 1765 1766

1767

1768

1769

1770 1771 1772

¹⁷³³ The following is our formal theorem.

1734 1735 1736 **Theorem 3.** Fix $\rho \in (0,1)$ and $a, b \in \mathbb{R}$ such that $\max\{|a|, |b|\} > 0$. Let $(A, \sigma^*) \sim \mathsf{ROS}_n(\rho, a, b)$. Let $y \sim \mathsf{BEC}(\sigma^*, \epsilon)$ or $y \sim \mathsf{BSC}(\sigma^*, \alpha)$, where

$$\lim_{n \to \infty} \frac{\log(1/\epsilon)}{\log n} = \beta \text{ and } \lim_{n \to \infty} \frac{\log(\frac{1-\alpha}{\alpha})}{\log n} = \beta$$

for some $\beta > 0$. Then I^* from (5) is well-defined and there is a degree-profiling algorithm (Algorithm 3) that returns the estimator $\hat{\sigma}_{dp}$ which achieves exact recovery whenever $I^* > 1$.

The following lemma plays a crucial role in the analysis of our degree profiling algorithm which formalizes the notion of receiving most of the labels correct.

Lemma F.4. Let $\sigma^* \in {\{\pm 1\}}^n$ be sampled such that each entry is i.i.d. with $\mathbb{P}(\sigma_i^* = +1) = \rho$. Condition on σ^* such that the event E from (9) holds. For any $\beta > 0$, we let $y \sim \mathsf{BEC}(\sigma^*, \epsilon_n)$ or $y \sim \mathsf{BSC}(\sigma^*, \alpha_n)$ for ϵ_n and α_n scales as described in Theorem 3. Define $S_+ = \{i : y_i = +1\}$ and $S_- = \{i : y_i = -1\}$. Then with probability 1 - o(1)

$$\max\{|C_{+} \setminus S_{+}|, |C_{-} \setminus S_{-}|\} = O\left(\frac{n}{\log^{10} n}\right)$$

1751 *Proof.* First, recall that conditioned on the even E about σ^* , we have $|C_+| = \Theta(n)$ and $|C_-| = \Theta(n)$. 1752 $\Theta(n)$. We now consider the two types of side information.

• BEC side information: Observe that

$$\mathbb{E}[|C_+ \setminus S_+|] = \sum_{i \in C_+} \mathbb{P}(y_i = 0) = |C_+|\epsilon_n = n^{-\beta}|C_+| \le n^{1-\beta}$$

Then Markov's inequality immediately implies that, with probability $1 - O(n^{-\beta/2})$, we have

$$|C_+ \setminus S_+| \le n^{1-\beta/2} = O(n/\log^{10} n)$$

Similarly, we also have $\mathbb{E}[|C_{-} \setminus S_{-}|] = n^{-\beta} |C_{-}| \le n^{1-\beta}$. Thus, applying Markov's inequality again implies $|C_{-} \setminus S_{-}| = O(n/\log^{10} n)$ with probability $1 - O(n^{-\beta/2})$. A simple union bound over these two events implies, with probability $1 - O(n^{-\beta/2}) = 1 - o(1)$

$$\max\{|C_{+} \setminus S_{+}|, |C_{-} \setminus S_{-}|\} = O\left(\frac{n}{\log^{10} n}\right)$$

• BSC side information: Under BSC side information, $\mathbb{E}[|C_+ \setminus S_+|] = \alpha_n |C_+| \le n^{1-\beta}$ and $\mathbb{E}[|C_- \setminus S_-|] = \alpha_n |C_-| \le n^{1-\beta}$. Therefore, using the Markov's inequality for both of these sets along with a union bound immediately implies that with probability $1 - O(n^{-\beta/2}) = 1 - o(1)$,

$$\max\left\{\left|C_{+}\setminus S_{+}\right|,\left|C_{-}\setminus S_{-}\right|\right\}=O\left(\frac{n}{\log^{10}n}\right).$$

1773 1774 1775

To bound the effect of the error terms when we make z^{dp} approximation to z^* , we need another technical lemma whose proof we include at the end of this section.

1778 Lemma F.5. Consider $\rho \in (0,1)$ and $a, b \in \mathbb{R}$ such that $\max\{|a|, |b|\} > 0$ Let $(A, \sigma^*) \sim \operatorname{ROS}_n(\rho, a, b)$. Condition on σ^* such that the event E holds. Fix any set $T \subset C_+$ or $T \subset C_-$ 1780 such that $|T| = O(n/\log^{10} n)$. Then for any $i \in [n]$, let us define $Y_i = \sqrt{\frac{\log n}{n}} \sum_{j \in T} A_{ij}$. 1781 $\mathbb{P}(\forall i \in [n] : |Y_i| \le 1) \ge 1 - O(n^{-3})$. Using this lemma, we now show that the degree profiling vector z^{dp} is a good approximation to the genie score vector z^* in ℓ_{∞} norm.

1785 Lemma F.6. Fix $\rho \in (0, 1)$ and a, b such that $\max\{|a|, |b|\} > 0$. Let $(A, \sigma^*) \sim \mathsf{ROS}_n(\rho, a, b)$ and **1786** condition on σ^* such that E from (9) holds. For $\beta > 0$, let $y \sim \mathsf{BEC}(\sigma^*, \epsilon_n)$ or $y \sim \mathsf{BSC}(\sigma^*, \alpha_n)$ **1787** where ϵ_n and α_n scales as described in Theorem 3. Let z^* and z^{dp} respectively be the genie score **1788** vector and the degree-profiling score vector produced by Algorithm 3 for the corresponding model **1789** of side information. Then (irrespective of the parameter values), with probability 1 - o(1),

$$||z^* - z^{dp}||_{\infty} = O(1).$$

Proof. We first start by observing, in the case of BEC side information z^{dp} is just formed by overrid-1797 ing the entries of z from step 1 of Algorithm 3 with $+\infty$ or $-\infty$ depending on the side information 1798 label being +1 or -1. Also, for BSC side information, $z^{dp} = z + \log\left(\frac{1-\alpha_n}{\alpha_n}\right) y$. By Lemmas 4.2, 1800 this is precisely how the genie score vector z^* in the respective model of side information relates to 1801 the genie score vector without side information which we denote by z'.

1802 Therefore, to show the lemma, it suffices to show that, with probability 1 - o(1),

$$||z'-z||_{\infty} = O(1).$$

• BEC side information:

$$\leq \max_{i \in [n]} \left| a(a-b)\sqrt{\frac{\log n}{n}} \sum_{j \in C_+ \setminus S_+} A_{ij} \right| + \left| b(a-b)\sqrt{\frac{\log n}{n}} \sum_{j \in C_- \setminus S_-} A_{ij} \right| + O(1), \tag{18}$$

where the last step follows from the triangle inequality. By Lemma F.4, both $|C_+ \setminus S_+|$ and $|C_- \setminus S_-|$ are bounded by $O(n/\log^{10} n)$ with probability 1 - o(1). Moreover, these sets are chosen only based on the side information y and hence independent of A, conditioned on σ^* . Using Lemma F.5 for these set $C_+ \setminus S_+$ and $C_- \setminus S_-$ as T, and using a union bound, we obtain that with probability 1 - o(1)

$$\max_{i \in [n]} \left| a(a-b)\sqrt{\frac{\log n}{n}} \sum_{j \in C_+ \setminus S_+} A_{ij} \right| = O(1) \quad \text{and} \quad \max_{i \in [n]} \left| b(a-b)\sqrt{\frac{\log n}{n}} \sum_{j \in C_- \setminus S_-} A_{ij} \right| = O(1).$$

Substituting these bounds in (18), with probability 1 - o(1)

 $\left\|z' - z\right\|_{\infty} = O(1).$

1836 • BSC side information: $||z' - z||_{\infty} = \max_{i \in [n]} |z'_i - z_i|$ 1838 1839 $= \max_{i \in [n]} \left| a(a-b) \sqrt{\frac{\log n}{n}} \sum_{i \in C \setminus S} A_{ij} - a(a-b) \sqrt{\frac{\log n}{n}} \sum_{i \in S \setminus \setminus C} A_{ij} \right|$ 1841 $+ b(a-b)\sqrt{\frac{\log n}{n}} \sum_{j \in C_{-} \setminus S_{-}} A_{ij} - b(a-b)\sqrt{\frac{\log n}{n}} \sum_{j \in S_{-} \setminus C_{-}} A_{ij} + \log\left(\frac{\rho}{1-\rho}\right) + o(1) \bigg|$ (substituting z' from Lemma F.1 and z from Algorithm 3) 1843 1845 $= \max_{i \in [n]} \left| (a-b)^2 \sqrt{\frac{\log n}{n}} \sum_{j \in C_+ \setminus S_+} A_{ij} - (a-b)^2 \sqrt{\frac{\log n}{n}} \sum_{j \in C_- \setminus S_-} A_{ij} + \log\left(\frac{\rho}{1-\rho}\right) + o(1) \right|$ 1849 (since $S_+ \setminus C_+ = C_- \setminus S_-$ and $S_- \setminus C_- = C_+ \setminus S_+$) $\leq \max_{i \in [n]} \left| (a-b)^2 \sqrt{\frac{\log n}{n}} \sum_{j \in C_+ \backslash S_+} A_{ij} \right| + \left| (a-b)^2 \sqrt{\frac{\log n}{n}} \sum_{i \in C_- \backslash S_-} A_{ij} \right| + O(1),$ 1851 (19)where in the last step, we used the triangle inequality. Again by similar arguments, first 1855 using Lemma F.4, both $|C_+ \setminus S_+|$ and $|C_- \setminus S_-|$ is $O(n/\log^{10} n)$ with probability 1 - o(1). Using Lemma F.5 for these set $C_+ \setminus S_+$ and $C_- \setminus S_-$ further implies that, with probability 1857 1 - o(1) $\max_{i \in [n]} \left| (a-b)^2 \sqrt{\frac{\log n}{n}} \sum_{i \in C_+ \setminus S_+} A_{ij} \right| = O(1) \quad \text{and} \quad \max_{i \in [n]} \left| (a-b)^2 \sqrt{\frac{\log n}{n}} \sum_{i \in C_- \setminus S_+} A_{ij} \right| = O(1).$ 1860 1862 Substituting these bounds in (19), with probability 1 - o(1)1864 $||z' - z||_{\infty} = O(1).$ 1866 1868 Finally, we prove Theorem 3. 1870 *Proof of Theorem 3.* We have already verified in Appendix C that I^* well-defined for BEC and BSC 1871 channels where ϵ and α scales as described in the theorem statements. When $\beta > 0$, by Lemma F.6 1872 that with probability 1 - o(1), 1873 $||z^* - z^{dp}||_{\infty} = O(1).$ 1874 1875 Above the IT threshold by Lemma D.1 and union bound over $i \in [n]$, there exists $\delta > 0$ such that 1876

$$\mathbb{P}\left(\min_{i\in[n]}\sigma_i^* z_i^* > \delta \log n\right) = 1 - o(1)$$

Taking a union bound, there exists $\varsigma > 0$ such that

1877 1878

1881 1882

1885

1889

$$\mathbb{P}\left(\min_{i\in[n]}\sigma_i^* z_i^{dp} > \varsigma \log n\right) = 1 - o(1)$$

Observing $\hat{\sigma}_{dp} = \text{sgn}(z^{dp})$, we obtain $\hat{\sigma}_{dp}$ achieves exact recovery, i.e.

$$\lim_{n \to \infty} \mathbb{P}\left(\hat{\sigma}_{dp} \text{ succeeds}\right) = 1.$$

We now return to the deferred proof.

Proof of Lemma F.5. First of all, observe that Y_i is a Gaussian random variable. Let us say $Y_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$, for some μ_i and $\sigma_i^2 > 0$. Applying Lemma B.2 for any Y_i (after renormalizing) yields

$$\mathbb{P}\left(\frac{|Y_i - \mu_i|}{\sigma_i} \ge 4\sqrt{\log n}\right) \le e^{-8\log n} = n^{-8}.$$

1895 Rearrangement of the terms using the triangle inequality along with a union bound over all $i \in [n]$ 1896 gives us

$$\mathbb{P}\left(\exists i \in [n] : |Y_i| \ge |\mu_i| + 4\sigma_i \sqrt{\log n}\right) \le n^{-7}$$

1899 Therefore it simply suffices to show that for every $i \in [n]$, we have $|\mu_i| + 4\sigma_i \sqrt{\log n} \le 1$. Indeed, 1900 we will show that these terms are o(1). To this end, first consider the term $4\sigma_i \sqrt{\log n}$ for any $i \in [n]$. 1901 Recall that Y_i is the sum of at most |T| i.i.d. Gaussian random variables, all with variance 1, scaled 1902 by $\sqrt{\log n/n}$. Therefore,

$$\sigma_i^2 \le \frac{\log n}{n} |T| = O\left(\frac{1}{\log^9 n}\right) \implies 4\sigma_i \sqrt{\log n} = O\left(\frac{1}{\log^2 n}\right) = o(1),$$

where we used $|T| = O(n/\log^{10} n)$. We now show that $|\mu_i| = o(1)$ too for all $i \in [n]$, which requires some casework. First consider $T \subset C_+$, then for any $i \in C_+$:

$$|\mu_i| \le \sqrt{\frac{\log n}{n}} |T| a^2 \sqrt{\frac{\log n}{n}} = O\left(\frac{1}{\log^9 n}\right),$$

1910 Similarly, for any $i \in C_-$:

$$\mu_i| = \sqrt{\frac{\log n}{n}} |T| |ab| \sqrt{\frac{\log n}{n}} = O\left(\frac{1}{\log^9 n}\right).$$

Exactly following the same arguments, we also get the same bounds on μ_i even when $T \subset C_-$. Overall, we established that

 $\mathbb{P}(\forall i \in [n] : |Y_i| \le 1) \ge 1 - O(n^{-7}).$

1916 1917

1894

1898

1903 1904

1909

1911 1912 1913

1918

1929 1930 1931

¹⁹¹⁹ G PROOFS AND ALGORITHMS FOR SBM.

We follow the same structure: in Appendix G.1, we derive the formula for genie scores when no side information is available. In Appendix G.2, we present our spectral algorithm with the optimality proof. Finally, in Appendix G.3, we do the degree profiling algorithm.

1924 1925 G.1 GENIE SCORES' FORMULA WHEN NO SIDE INFORMATION

1926 We begin by showing that, with high probability, all the vertices have degrees logarithmic in *n*.

1927 Lemma G.1. Let $\rho \in (0, 1)$ and $a_1, a_2, b > 0$. Let $(A, \sigma^*) \sim \text{SBM}_n(\rho, a_1, a_2, b)$. Condition on σ^* **1928** such that the event E from (20) holds then. For $c = 6 \max\{1, a_1, a_2, b\}$ let

$$E_1 = \left\{ \forall i : \sum_{j \in [n]} A_{ij} \le c \log n \right\};$$
(20)

1932 1933 then with $\mathbb{P}(E_1) = 1 - O(n^{-3})$.

Proof. Note that the entries of A (up to symmetry) are independent conditioned on σ^* . Therefore, for any $i \in [n]$, the *i*th row has independent Bernoulli entries with means either p_1, p_2 or q, where $(p_1, p_2, q) = (a_1, a_2, b) \log n/n$. Therefore, defining $X \sim \text{Binom}(n, \tau \log n/n)$, where $\tau = \max\{a_1, a_2, b\}$, we have that X stochastically dominates $\sum_{j \in [n]} A_{ij}$, for any $i \in [n]$. Then applying the Chernoff bound for Binomial random variables (Mitzenmacher & Upfal, 2017, Theorem 4.4, Equation 4.3) we get, for any $i \in [n]$

1943

$$\mathbb{P}\left(\sum_{j\in[n]} A_{ij} > 6\max\{1,\tau\}\log n\right) \le \mathbb{P}\left(X > 6\max\{1,\tau\}\log n\right) \le 2^{-6\log n} = O(n^{-4})$$

Taking a union bound over all $i \in [n]$ yields the desired claim.

We next analyze the form of genie scores without side information.

Lemma G.2. Let $\rho \in (0,1)$ and $a_1, a_2, b > 0$. Let $(A, \sigma^*) \sim \text{SBM}_n(\rho, a_1, a_2, b)$. Denote $(p_1, p_2, q) := (a_1, a_2, b) \log n/n$. Then for any $i \in [n]$, the genie score can be written as

$$z_i^* = \log\left(\frac{p_1(1-q)}{q(1-p_1)}\right) \sum_{j \in C_+^{-i}} A_{ij} + \log\left(\frac{q(1-p_2)}{p_2(1-q)}\right) \sum_{j \in C_-^{-i}} A_{ij} + \log\left(\frac{\rho}{1-\rho}\right)$$

$$+ |C_{+}^{-i}| \log\left(\frac{1-p_{1}}{1-q}\right) + |C_{-}^{-i}| \log\left(\frac{1-q}{1-p_{2}}\right).$$

1953 Moreover, conditioned on the event E from (9) and E_1 from (20),

$$|z^* - Aw - \gamma \mathbf{1}_n||_{\infty} = O(1),$$

1956 where $w \in \mathbb{R}^n$ is a vector with entries $(w_+, w_-) := (\log(a_1/b), \log(b/a_2))$ on locations of C_+ and 1957 C_- respectively and $\gamma := (\rho(b-a_1) + (1-\rho)(a_2-b)) \log n$.

Proof. First of all, note that the SBM is a special case of the GBM model with $\mathcal{P}_+ \equiv \text{Bern}(p_1)$, $\mathcal{P}_- \equiv \text{Bern}(p_2)$ and $\mathcal{Q} \equiv \text{Bern}(q)$. Using Lemma 4.2 for this special case, for any $i \in [n]$

$$z_{i}^{*} = \log\left(\frac{\rho}{1-\rho}\right) + \sum_{i \in C_{+}^{-i}} \log\left(\frac{p_{1}^{A_{ij}}(1-p_{1})^{(1-A_{ij})}}{q^{A_{ij}}(1-q)^{(1-A_{ij})}}\right) + \sum_{j \in C_{-}^{-i}} \log\left(\frac{q^{A_{ij}}(1-q)^{(1-A_{ij})}}{p_{2}^{A_{ij}}(1-p_{2})^{(1-A_{ij})}}\right)$$
$$= \log\left(\frac{p_{1}(1-q)}{q(1-p_{1})}\right) \sum_{j \in C_{+}^{-i}} A_{ij} + \log\left(\frac{q(1-p_{2})}{p_{2}(1-q)}\right) \sum_{j \in C_{-}^{-i}} A_{ij} + \log\left(\frac{\rho}{1-\rho}\right)$$

$$+ |C_{+}^{-i}| \log\left(\frac{1-p_{1}}{1-q}\right) + |C_{-}^{-i}| \left(\frac{1-q}{1-p_{2}}\right).$$
(21)

To show the second part of the lemma, we further simplify

$$\left|\log\left(\frac{1-q}{1-p_1}\right)\right| = \left|\log\left(1+\frac{p_1-q}{(1-p_1)}\right)\right| = \left|\log\left(1+\frac{(a_1-b)\log n}{(1-p_1)n}\right)\right| = O\left(\frac{\log n}{n}\right).$$
 (22)
In the last inequality, we used $\frac{x}{1-q} < \log(1+q) < x$ for $q > -1$. Similarly

1973 In the last inequality, we used $\frac{x}{x+1} \le \log(1+x) \le x$ for x > -1. Similarly,

$$\left|\log\left(\frac{1-p_2}{1-q}\right)\right| = \left|\log\left(1+\frac{q-p_2}{(1-q)}\right)\right| = \left|\log\left(1+\frac{(b-a_2)\log n}{(1-q)n}\right)\right| = O\left(\frac{\log n}{n}\right).$$
 (23)

Recall the definition of event E from (9) and E_1 from (20). Conditioned on $E \cap E_1$, we simplify (21) using (22) and (23).

$$\log\left(\frac{p_{1}(1-q)}{q(1-p_{1})}\right) \sum_{j \in C_{+}^{-i}} A_{ij} + \log\left(\frac{q(1-p_{2})}{p_{2}(1-q)}\right) \sum_{j \in C_{-}^{-i}} A_{ij}$$

$$= \log\left(\frac{a_1}{b}\right) \sum_{j \in C_+^{-i}} A_{ij} + \log\left(\frac{b}{a_2}\right) \sum_{j \in C_-^{-i}} A_{ij} + O\left(\frac{\log n}{n}\right) \sum_{j \in [n]} A_{ij}$$

$$= \log\left(\frac{a_1}{b}\right) \sum_{j \in C_+^{-i}} A_{ij} + \log\left(\frac{b}{a_2}\right) \sum_{j \in C_-^{-i}} A_{ij} + o(1),$$
(24)

where the last equality followed by conditioning on E_1 . We also simplify

$$|C_{+}^{-i}|\log\left(\frac{1-p_{1}}{1-q}\right) = |C_{+}^{-i}|\log\left(1+\frac{q-p_{1}}{(1-q)}\right) = |C_{+}^{-i}|\log\left(1+\frac{(b-a_{1})\log n}{(1-q)n}\right)$$
$$= |C_{+}^{-i}|\left(\frac{(b-a_{1})\log n}{(1-q)n} + O\left(\frac{\log^{2} n}{n^{2}}\right)\right)$$

(using a Taylor expansion of log(1+x))

(25)

1995
1996
$$= (1 + O(n^{-1/3}))\rho n \left(\frac{(b-a_1)\log n}{(1-q)n} + O\left(\frac{\log^2 n}{n^2}\right)\right)$$

1997
$$= \rho(b - a_1) \log n + o(1)$$

Similarly,

$$C_{-i}^{-i} \log\left(\frac{1-q}{1-p_2}\right) = |C_{-i}^{-i}| \log\left(1 + \frac{p_2 - q}{(1-p_2)}\right) = (1-\rho)(a_2 - b) \log n + o(1)$$
(26)

Substituting (24), (25) and (26) into (21),

$$z_i^* = \log\left(\frac{a_1}{b}\right) \sum_{j \in C_+^{-i}} A_{ij} + \log\left(\frac{b}{a_2}\right) \sum_{j \in C_-^{-i}} A_{ij} + \rho(b - a_1) + (1 - \rho)(a_2 - b)\log n + O(1)$$
$$= w_+ \sum_{j \in C_+^{-i}} A_{ij} + w_- \sum_{j \in C_-^{-i}} A_{ij} + \gamma + O(1) = A_i \cdot w + \gamma + O(1).$$

Writing the above for all $i \in [n]$ in a vector notation, we obtain

$$\|z^* - (Aw + \gamma \mathbf{1}_n)\|_{\infty} = O(1). \quad \Box$$

2015 G.2 SPECTRAL ALGORITHM AND PROOF OF THEOREM 2

In this section, we present our spectral algorithm for the SBM that can emulate the genie. As discussed in Section 5, this requires taking an appropriate linear combination of eigenvectors such that the top two eigenvectors such that $c_1u_1^* + c_2u_2^*$ approximates w in the ℓ_{∞} norm. The vector wis a block vector with entries (w_+, w_-) on the locations of C_+ and C_- . Recall by Lemma G.2 that

$$w_{+} = \log\left(\frac{a_{1}}{b}\right)$$
 and $w_{-} = \log\left(\frac{b}{a_{2}}\right)$

We first present a subroutine that finds these coefficients (c_1, c_2) . We will introduce the formal correctness of the subroutine later, but first we provide informal discussion as to how these coefficients are computed. Roughly speaking, both u_1^* and u_2^* also have a block structure, and therefore, finding (c_1, c_2) just corresponds to solving a system of 2×2 linear equations. Also, the coefficients do not depend on the locations of σ^* with +1 or -1 labels, so exploiting this fact we just do calculation as if C_+ is on the first $\lfloor \rho n \rfloor$ vertices and compute the proxy for actual A^* . This results into the following subroutine.

Algorithm 4 Find Linear Combination Coefficients

Input: The parameter set (ρ, a_1, a_2, b) such that $a_1 a_2 \neq b^2$ (Rank-2) and the graph size n.

Output: The desired linear combination (c_1, c_2) .

1: Let $S \subseteq [n]$ such that $S = \{i : i \leq \rho n\}$ and compute the matrix $B \in \mathbb{R}^{n \times n}$ such that

$$B_{ij} = \begin{cases} a_1 \log n/n & \text{if } i, j \in S; \\ b \log n/n & \text{if } i \in S \text{ but } j \notin S; \\ a_2 \log n/n & \text{if } i \notin S \text{ and } j \notin S. \end{cases}$$

2: Compute the two eigenpairs $(\tilde{\lambda}_1, \tilde{v}_1)$ and $(\tilde{\lambda}_2, \tilde{v}_2)$ of B (note that B is rank-2).

3: Let $w \in \mathbb{R}^n$ be the block vector such that:

$$w_i = \begin{cases} w_+ = \log(a_1/b), & \text{if } i \in S; \\ w_- = \log(b/a_2) & \text{if } i \notin S. \end{cases}$$

2046 4: Return $(c_1, c_2) \in \mathbb{R}^2$ that satisfies

$$c_1\left(\frac{\tilde{v}_1}{\tilde{\lambda}_1}\right) + c_2\left(\frac{\tilde{v}_2}{\tilde{\lambda}_2}\right) = w.$$
(27)

Both \tilde{v}_1 and \tilde{v}_2 are block-vectors and they are linearly independent. Thus, Finding (c_1, c_2) corresponds to solving a system of 2×2 linear equations.

2052 We note that in the rank-2 case when $a_1a_2 \neq b^2$, the vectors \tilde{v}_1 and \tilde{v}_2 have block structure and 2053 are linearly independent. Therefore, it is possible to span any vector with block structure, and in 2054 particular, even w. We now propose our spectral algorithm in the rank-2 case. 2055 Algorithm 5 Spectral recovery algorithm for SBM (Rank-2) 2056 2057 **Input:** An $n \times n$ observation matrix A and parameters (ρ, a_1, a_2, b) such that $a_1 a_2 \neq b^2$. Optionally, 2058 side information $y \in \mathcal{Y}^n$ such that we can compute the likelihood laws \mathcal{S}_+ and \mathcal{S}_- . 2059 **Output:** An estimate of community assignments $\hat{\sigma}_{spec}$. 2060 1: Compute leading eigenpairs. Compute the top eigenpair of A, denoted by (λ_1, u_1) , where 2061 $|\lambda_1| \geq \cdots \geq |\lambda_n|.$ 2062 2: Compute coefficients of linear combination. Run Algorithm 4 to find (c_1, c_2) and 2063 2064 $\gamma := (\rho(b - a_1) + (1 - \rho)(a_2 - b)) \log n.$ 2065 3: Compute spectral scores. For any $s = (s_1, s_2) \in \{\pm 1\}^2$, prepare the spectral score vectors 2066 • No side information: 2067 $z^{(s)} = s_1 c_1 u_1 + s_2 c_2 u_2 + \gamma \mathbf{1}_n.$ • Side information: 2069 $z^{(s)} = s_1 c_1 u_1 + s_2 c_2 u_2 + \gamma \mathbf{1}_n + \log\left(\frac{S_+}{S}(y)\right).$ 2071 4: Remove sign ambiguity. For each $s \in \{\pm 1\}^2$, let $\hat{\sigma}^{(s)} = \operatorname{sgn}(z^{(s)})$. 2073 • No side information: Return $\hat{\sigma}_{\text{spec}} = \arg \max_{\{\hat{\sigma}^{(s)}:s \in \{\pm 1\}\}} \mathbb{P}(\sigma^* = \hat{\sigma}^{(s)} \mid A).$ 2074 2075 • BEC or BSC side information: Return $\hat{\sigma}_{\text{spec}} = \arg \max_{\{\hat{\sigma}^{(s)}:s \in \{\pm 1\}\}} \mathbb{P}(\sigma^* = \hat{\sigma}^{(s)} \mid A, y).$ 2076 2077 Finally, when $a_1a_2 = b^2$, from Lemma G.2 we have $w_+ = w_- = \log\left(\frac{a_1}{b}\right)$. In this case, we can 2078 just use the deterministic vector along $\mathbf{1}_n$ to emulate the genie-score vector. Strictly speaking, in 2079 this degenerate case, we do not even need a spectral strategy to achieve optimality. Despite this, we 2080 just refer to this as a spectral algorithm in the rest of the analysis for simplicity of exposition. 2081 Algorithm 6 (Spectral) Recovery algorithm for SBM (Rank-1) 2082 **Input:** An $n \times n$ observation matrix A and parameters (ρ, a_1, a_2, b) such that $a_1a_2 = b^2$. Optionally, 2083 side information $y \in \mathcal{Y}^n$ such that we can compute the likelihood laws \mathcal{S}_+ and \mathcal{S}_- .. 2084 2085 **Output:** An estimate of community assignments $\hat{\sigma}_{\text{spec}}$. 2086 1: Let $c := \log\left(\frac{a_1}{b}\right), \quad \gamma := (\rho(b - a_1) + (1 - \rho)(a_2 - b))\log n.$ 2089 2: Prepare the spectral score vector z^{spec} as follows. 2090 • No side information: 2091 $z^{\text{spec}} = cA\mathbf{1}_n + \gamma \mathbf{1}_n.$ 2092 • Side information: $z^{\text{spec}} = cA\mathbf{1}_n + \gamma \mathbf{1}_n + \log\left(\frac{S_+}{S_-}(y)\right).$ 2093 2094 2095 3: Return $\hat{\sigma}_{\text{spec}} = \text{sgn}(z^{\text{spec}})$. 2096 2097 We now show that the score vectors formed by these algorithms approximate the genie score vector 2098 z^* in each case. **Lemma G.3.** Consider $\rho \in (0,1)$ and $a_1, a_2, b > 0$. Let $(A, \sigma^*) \sim \mathsf{SBM}_n(\rho, a_1, a_2, b)$ and 2100 condition on σ^* such that E from (9) holds. Optionally, let $y \sim Sl(\sigma^*, S_+, S_-)$ for the channel 2101 laws $(\mathcal{S}_+, \mathcal{S}_-)$. Let z^* be the genie score vector for the corresponding model. Then with probability 2102 1 - o(1)

2103
2104
• Case
$$a_1a_2 \neq b^2$$
: for some $s = (s_1, s_2) \in \{\pm 1\}^2$
2105
 $\left\| z^* - z^{(s)} \right\|_{\infty} = o(\log n).$

• Case
$$a_1a_2 = b^2$$
:

Proof. Recall from Lemma 4.2, how the genie score changes in the presence of side information. In Algorithm 5 (step 3) and Algorithm 6 (step 2), this is precisely how the spectral score vectors are updated from the case when no side information is present. Therefore, it suffices to show that the score approximation holds in the case when no side information is present, which now will be the focus of the proof. The argument is exactly analogous to the one done used in Lemma F.3. We now discuss the rank 1 and rank 2 cases one by one.

 $\|z^* - z^{\operatorname{spec}}\|_{\infty} = o(\log n)$

Rank-1 case: $a_1a_2 = b^2$. In this case, when no side information is present

$$||z^* - z^{\text{spec}}||_{\infty} = ||z^* - \log(a_1/b)A\mathbf{1}_n - \gamma\mathbf{1}_n||_{\infty} = O(1),$$

where the last equation follows from Lemma G.2 and using that $\frac{a_1}{b} = \frac{b}{a_2}$.

Rank-2 case: $a_1a_2 \neq b^2$. Fix $(s_1, s_2) \in \{\pm 1\}^n$ to be the signs for which the high probability event in Lemma E.4 holds. Define w to be the vector from Lemma G.2 with entries $(w_+, w_-) =$ $(\log(a_1/b), \log(b/a_2))$ on the locations of C_+ and C_- respectively. In this case, using Lemma E.4, we will show that with probability 1 - o(1), we have

$$\|Aw - s_1c_1u_1 - s_2c_2u_2\|_{\infty} = o(\log n).$$
⁽²⁸⁾

Combing this (28) along with Lemma G.2 implies the desired result: with probability 1 - o(1),

$$\begin{aligned} \left\| z^* - z^{(s)} \right\|_{\infty} &\leq \left\| z^* - Aw - \gamma \mathbf{1}_n \right\|_{\infty} + \left\| Aw + \gamma \mathbf{1}_n - z^{(s)} \right\|_{\infty} \\ &= \left\| z^* - Aw - \gamma \mathbf{1}_n \right\|_{\infty} + \left\| Aw - s_1 c_1 u_1 - s_2 c_2 u_2 \right\|_{\infty} \quad \text{(step 3 of Algorithm 5)} \\ &= o(\log n). \end{aligned}$$

It remains to show (28). Note that, we calculate (c_1, c_2) in Algorithm 4 using the matrix B where community sizes are exactly ρn . But, condition on E, we have community sizes $(1 + o(1))\rho n$ and $(1 + o(1))(1 - \rho)n$. Also, in A^{*}, we have zero diagonal, where as, the matrix B has diagonal entries of the order of $O(\log n/n)$. These changes only affect the eigenvalues by the multiplicative factor of (1 + o(1)) by Weyl's inequality. Therefore,

$$\lambda_1^* = (1 + o(1))\tilde{\lambda}_1 \text{ and } \lambda_2^* = (1 + o(1))\tilde{\lambda}_2.$$

Moreover, the entries of u_1^* in location of C_+ are (1 + o(1)) factor of the entries of \tilde{v}_1 in the location S. By the same argument, one can say the same for u_1^* in locations of C_- and \tilde{v}_1 in locations of $[n] \setminus S$, and also for u_2^* and \tilde{v}_2 . From the way we calculate (c_1, c_2) in (27), we have

$$c_1\left(\frac{\tilde{v}_1}{\tilde{\lambda}_1}\right) + c_2\left(\frac{\tilde{v}_2}{\tilde{\lambda}_2}\right) = w$$

Then the above discussion implies

2144
2145
$$w_{+} = (1+o(1))\left(\frac{c_{1}u_{1,i}^{*}}{\lambda_{1}^{*}} + \frac{c_{2}u_{2,i}^{*}}{\lambda_{2}^{*}}\right), \text{ for } i \in C_{+} \text{ and } w_{-} = (1+o(1))\left(\frac{c_{1}u_{1,i}^{*}}{\lambda_{1}^{*}} + \frac{c_{2}u_{2,i}^{*}}{\lambda_{2}^{*}}\right), \text{ for } i \in C_{-}$$
2146

Finally using Lemma E.4, with probability 1 - o(1)

$$\begin{aligned} & \|Aw - s_1c_1u_2 - s_2c_2u_2\|_{\infty} \le \left\|Aw - c_1\frac{Au_1^*}{\lambda_1^*} - c_2\frac{Au_2^*}{\lambda_2^*}\right\|_{\infty} + o(1) \\ & = \left\|(1 + o(1))\left(c_1\frac{Au_1^*}{\lambda_1^*} + c_2\frac{Au_2^*}{\lambda_2^*}\right) - c_1\frac{Au_1^*}{\lambda_1^*} - c_2\frac{Au_2^*}{\lambda_2^*}\right\|_{\infty} + o(1) \\ & = o(1)\left(\max\|A_{i,i}\|_{1}\right) \cdot \left(\left\|\frac{c_1u_1^*}{u_1^*}\right\|_{\infty} \lor \left\|\frac{c_2u_2^*}{u_2^*}\right\|_{\infty}\right) \end{aligned}$$

$$= o(1) \left(\max_{i \in [n]} \|A_{i\cdot}\|_1 \right) \cdot \left(\left\| \frac{1 + 1}{\lambda_1^*} \right\|_{\infty} \vee \left\| \frac{1 + 2 + 2}{\lambda_2^*} \right\|_{\infty} \right)$$

where the last step follows by using Lemma G.1 and (c_1, c_2) are chosen in such a way that the entries of vectors $\frac{c_1u_1^*}{\lambda_1^*}$ and $\frac{c_2u_2^*}{\lambda_2^*}$ are O(1) by (27).

We finally combine all the pieces and give the proof of Theorem 2.

 $= o(\log n),$

2160	Proof of Theorem 2. We break into cases.
2162	1. Rank 1, i.e. $a_1a_2 = b^2$: First note that Algorithm 6 creates a score vector $z^{\text{spec}} \in \mathbb{R}^n$ such
2163	that $ z^* - z^{\text{spec}} _{\infty} = o(\log n)$ by Lemma G.3. Using Lemma D.1, whenever $I^* >$, there
2164	exists $\varsigma > 0$ such that
2165	$\mathbb{P}(\sigma_i^* z_i^{prec} \le \varsigma \log n) = o(n^{-1}).$
2166	Taking a union bound,
2167	$\mathbb{P}(\forall i \in [n], \sigma_i^* z_i^{\text{spec}} > \varsigma \log n) = 1 - o(1).$
2168	Eight the electric function \hat{c} $constants \hat{c}$ $constants (spec) (step 2) this immediately implies$
2169	Finally, the algorithm outputs $\sigma_{\text{spec}} = \text{sgn}(z^{+r+1})$ (step 5), this immediately implies
2170	$\mathbb{P}(\hat{\sigma}_{ ext{spec}} = \sigma^*) = 1 - o(1).$
2172	2. Doub 2 is a $\frac{1}{2}$ We first note that whenever $I^* > 1$ the estimator \hat{c} - solving
2173	2. Rank 2, i.e. $a_1 a_2 \neq 0^-$: we first note that whenever $T^* > 1$, the estimator σ_{MAP} achieves exact recovery by Proposition 3.1. That is with high probability, we have $\hat{\sigma}_{MAP} = \sigma^*$ unless
2174	$a_1 = a_2$ and no side information, in which case $\hat{\sigma}_{MAP} \in \{\pm \sigma^*\}$. Recall that Algorithm 5
2175	creates four candidates for σ^* and chooses the one with maximum posterior probability.
2176	Due to statistical achievability, it suffices to show that one of the $\{\sigma^{(s)}: s = (s_1, s_2) \in$
2177	$\{\pm 1\}^2\}$ maintained by the algorithm is σ^* with high probability. To this end, recall Lemma
2178	G.3 that with probability $1 - o(1)$,
2179	$\min \ z^* - z^{(s)}\ = o(\log n).$
2180	$s \in \{\pm 1\}$ $\ \sim \rangle \sim \ _{\infty}$
2101	Combining this with Lemma D.1, there exists $s^* \in \{\pm 1\}^2$ and $\varsigma > 0$ such that
2183	
2184	$\mathbb{P}\left(\min \sigma_i^* z_i^{(s^*)} > \varsigma \log n\right) = 1 - o(1),$
2185	$i \in [n]$
2186	after taking a union bound. As $\hat{\sigma}^{(s^*)} = \operatorname{sgn}(z^{(s^*)})$ in step 4, we obtain $\hat{\sigma}^{(s^*)} = \sigma^*$ with
2187	high probability. Overall, we established that $\hat{\sigma}_{spec}$ achieves exact recovery above the IT
2188	threshold.
2189	
2190	
2191	G.3 DEGREE-PROFILING ALGORITHM FOR BEC AND BSC CHANNELS
2193	Algorithm 7 Degree-Profiling algorithm for SBM in the presence of BEC or BSC side information.
2194 2195	Input: An $n \times n$ observation matrix A and parameters (ρ, a_1, a_2, b) . The BEC side information y with parameter ϵ or BSC side information y with parameter α .
2196	Output: An estimate of community assignments $\hat{\sigma}_{dp}$.
2197	1: Let $S_1 := \{i : y_i = +1\}$ $S_1 := \{i : y_i = -1\}$ and $\gamma := (o(b - a_1) + (1 - o)(a_2 - b))\log n$
2198 2199	Compute $z \in \mathbb{R}^n$ such that, for every $i \in [n]$
2200 2201 2202	$z_i = \log\left(\frac{a_1}{b}\right) \sum_{j \in S_+} A_{ij} + \log\left(\frac{b}{a_2}\right) \sum_{j \in S} A_{ij} + \gamma.$
2203 2204 2205	 2: Prepare the degree-profile score vector z^{dp} as follows. BEC side information: For any i ∈ [n],
2206	$\int z_i$ if $u_i = 0$:
2207	$z_{i}^{dp} = \begin{cases} z_{i}^{m} & z_{i}^{m} & z_{i}^{m} \\ +\infty & \text{if } y_{i} = +1 \end{cases}$
2208	$\int_{-\infty}^{+\infty} if y_i = -1;$
2209	$y_i = y_i$

• BSC side information:

$$z^{\rm dp} = z + \log\left(\frac{1 - \alpha_n}{\alpha_n}\right) y$$

2213 3: Return $\hat{\sigma}_{dp} = \operatorname{sgn}(z^{dp})$.

2210 2211

The success of the degree-profiling algorithm is formalized in the following lemma.

Theorem 4. Let $\rho \in (0,1)$ and $a_1, a_2, b > 0$. Let $(A, \sigma^*) \sim \mathsf{SBM}_n(\rho, a_1, a_2, b)$. Let $y \sim \mathsf{BEC}(\sigma^*, \epsilon)$ or $y \sim \mathsf{BSC}(\sigma^*, \alpha)$, where

2218 2219

2220

2230 2231 2232

2240 2241 2242

2249 2250 2251

2252

2253

2254 2255

2262

2263

2265 2266

$$\lim_{n \to \infty} \frac{\log(1/\epsilon)}{\log n} = \beta \text{ and } \lim_{n \to \infty} \frac{\log(\frac{1-\alpha}{\alpha})}{\log n} = \beta$$

for some $\beta > 0$. Then I^* from (5) is well-defined and there is a degree-profiling algorithm (Algorithm 7) that returns the estimator $\hat{\sigma}_{dp}$ which achieves exact recovery whenever $I^* > 1$.

Again to prove that our approximation is good enough when using S_+ and S_- as proxies for the actual communities, we will need a technical lemma to bound the error terms, whose proof we include the at the end of this section.

Lemma G.4. Consider $\rho \in (0, 1)$ and $a_1, a_2, b > 0$. Let $(A, \sigma^*) \sim \mathsf{SBM}_n(\rho, a_1, a_2, b)$. Condition on σ^* such that the event E from (9) holds. Fix any set $T \subset C_+$ or $T \subset C_-$ such that $|T| = O(n/\log^{10} n)$. Then for any $i \in [n]$, define $Y_i := \sum_{j \in T} A_{ij}$. Then

$$\mathbb{P}\left(\forall i \in [n] : Y_i \le \frac{\log n}{\log \log n}\right) \ge 1 - O(n^{-3}).$$

Using this, we now show that the degree-profiling scores are indeed good approximations of the actual genie scores.

Lemma G.5. Fix $\rho \in (0,1)$ and $a_1, a_2, b > 0$. Let $(A, \sigma^*) \sim \text{SBM}_n(\rho, a_1, a_2, b)$ and condition on σ^* such that E from (9) holds. For $\beta > 0$, let $y \sim \text{BEC}(\sigma^*, \epsilon_n)$ or $y \sim \text{BSC}(\sigma^*, \alpha_n)$ where ϵ_n and α_n scales as described in Theorem 4. Let z^* and z^{dp} respectively be the genie score vector and the degree-profiling score vector produced by Algorithm 7 for the corresponding model of side information. Then with probability 1 - o(1),

$$\left\|z^* - z^{\mathsf{dp}}\right\|_{\infty} = O\left(\frac{\log n}{\log \log n}\right)$$

2243 *Proof.* The proof has similar calculations as in Lemma F.6 for ROS. We again first start by noting 2244 that z^{dp} is just formed by overriding the entries of z from step 1 of Algorithm 7 with $+\infty$ or $-\infty$ 2245 depending on the side information label being +1 or -1 under BEC channel. Also, in step 2 under 2246 the BSC channel, we have $z^{dp} = z + \log\left(\frac{1-\alpha_n}{\alpha_n}\right) y$. Recall Lemma 4.2 as to how the genie scores 2247 change under both types of side information. It suffices to show that, with probability 1 - o(1),

$$\left\|z' - z\right\|_{\infty} = O\left(\frac{\log n}{\log \log n}\right)$$

where z' is the genie score vector without side information and z is formed in step 1. Recall that E_1 holds with probability $1 - O(n^{-3})$ from (20). Using Lemma G.2 along with the triangle inequality we obtain the following.

• BEC side information:

$$\begin{aligned} \|z' - z\|_{\infty} &= \max_{i \in [n]} |z'_i - z_i| \\ &\leq \max_{i \in [n]} \left| \log \left(\frac{a_1}{b}\right) \sum_{j \in C_+ \setminus S_+} A_{ij} + \log \left(\frac{b}{a_2}\right) \sum_{j \in C_- \setminus S_-} A_{ij} \right| + O(1) \end{aligned}$$

where we substitute z' from step 1. From Lemma F.4, both $|C_+ \setminus S_+|$ and $|C_- \setminus S_-|$ are bounded by $O(n/\log^{10} n)$ with probability 1 - o(1). These sets are independent of A and are chosen based on y. Using Lemma G.4 for these sets, with probability 1 - o(1)

$$\|z'-z\|_{\infty} \le \max_{i \in [n]} \left| \log\left(\frac{a_1}{b}\right) \sum_{j \in C_+ \setminus S_+} A_{ij} \right| + \left| \log\left(\frac{b}{a_2}\right) \sum_{j \in C_- \setminus S_-} A_{ij} \right| + O(1) = O\left(\frac{\log n}{\log \log n}\right).$$

• <u>BSC side information</u>: There are additional error terms caused by the sets $S_+ \setminus C_+$ and $S_{-} \setminus C_{-}$:

$$\|z'-z\|_{\infty} = \max_{i \in [n]} |z'_i - z_i| = \max_{i \in [n]} \left| \log\left(\frac{a_1}{b}\right) \sum_{j \in C_+ \setminus S_+} A_{ij} - \log\left(\frac{a_1}{b}\right) \sum_{j \in S_+ \setminus C_+} A_{ij} \right|$$

$$+ \log\left(\frac{b}{a_2}\right) \sum_{j \in C_- \setminus S_-} A_{ij} - \log\left(\frac{b}{a_2}\right) \sum_{j \in S_- \setminus C_-} A_{ij} \left| + O(1) \right|$$

T

(substituting z' from Lemma G.2 and z from step 1)

$$= O(1) \cdot \max_{i \in [n]} \left| \sum_{j \in C_+ \setminus S_+} A_{ij} + \sum_{j \in C_- \setminus S_-} A_{ij} \right| + O(1).$$

(since $S_+ \setminus C_+ = C_- \setminus S_-$ and $S_- \setminus C_- = C_+ \setminus S_+$)

Exactly the same argument of using Lemma F.4, both $|C_+ \setminus S_+|$ and $|C_- \setminus S_-|$ is $O(n/\log^{10} n)$ with probability 1 - o(1). Using Lemma G.4 for these sets, with probability 1 - o(1)

$$\left\|z' - z\right\|_{\infty} = O\left(\frac{\log n}{\log \log n}\right)$$

 $\log n$

We finally prove Theorem 4.

Proof of Theorem 4. We already discussed in Appendix C that I^* is well-defined. When $\beta > 0$, by Lemma G.5, with probability 1 - o(1),

$$\left\|z^* - z^{\mathrm{dp}}\right\|_{\infty} = O(1)$$

Above the IT threshold by Lemma D.1 and union bound over $i \in [n]$, there exists $\delta > 0$ such that

$$\mathbb{P}\left(\min_{i\in[n]}\sigma_i^* z_i^* > \delta \log n\right) = 1 - o(1)$$

Taking a union bound of these two events, there exists $\varsigma > 0$ such that

$$\mathbb{P}\left(\min_{i\in[n]}\sigma_i^* z_i^{\mathrm{dp}} > \varsigma \log n\right) = 1 - o(1).$$

Observing $\hat{\sigma}_{dp} = \text{sgn}(z^{dp})$ in step 3 of Algorithm 7, we obtain $\hat{\sigma}_{dp}$ achieves exact recovery, i.e.

$$\lim_{n \to \infty} \mathbb{P}\left(\hat{\sigma}_{dp} \text{ succeeds}\right) = 1.$$

We finally return to the skipped proof of a technical lemma.

Proof of Lemma G.4. Fix any set T according to the lemma and define $\{Y_i : i \in [n]\}$. Let $\tau = \max\{a_1, a_2, b\}$ and $Y \sim \operatorname{Binom}(|T|, \frac{\tau \log n}{n})$. Observe that for any $i \in [n]$, we have Y_i is stochastically dominated by Y. Applying the following Chernoff bound for Binomial random variable (Mitzenmacher & Upfal, 2017, Thereom 4.4): for any r > 1 and $X \sim \text{Binom}(N, p)$, $\mathbb{P}(X \ge rnp) \le (e/r)^{rnp}$, we obtain for any $i \in [n]$

$$\mathbb{P}\left(Y_i \le \frac{\log n}{\log\log n}\right) \le \mathbb{P}\left(Y \le \frac{\log n}{\log\log n}\right) \le \left(\frac{e|T|\tau \log n/n}{\log \log n}\right)^{\frac{\log n}{\log\log n}}$$

2315
2316
$$= O\left(\frac{\log\log n}{\log^{10} n}\right)^{\frac{\log n}{\log\log n}} = n^{-10+o(1)}.$$
2317

A simple union bound over all $i \in [n]$ gives us

2319
2320
$$\mathbb{P}\left(\forall i \in [n] : Y_i \le \frac{\log n}{\log \log n}\right) \ge 1 - O(n^{-3}). \quad \Box$$