

# PARTIAL RECOVERY FOR TOP- $k$ RANKING: OPTIMALITY OF MLE AND SUBOPTIMALITY OF THE SPECTRAL METHOD

BY PINHAN CHEN<sup>1,a</sup>, CHAO GAO<sup>1,b</sup> AND ANDERSON Y. ZHANG<sup>2,c</sup>

<sup>1</sup>Department of Statistics, University of Chicago, [pinhanchen@uchicago.edu](mailto:pinhanchen@uchicago.edu), [chaogao@uchicago.edu](mailto:chaogao@uchicago.edu)

<sup>2</sup>Department of Statistics, University of Pennsylvania, [ayz@wharton.upenn.edu](mailto:ayz@wharton.upenn.edu)

Given partially observed pairwise comparison data generated by the Bradley–Terry–Luce (BTL) model, we study the problem of top- $k$  ranking. That is, to optimally identify the set of top- $k$  players. We derive the minimax rate with respect to a normalized Hamming loss. This provides the first result in the literature that characterizes the partial recovery error in terms of the proportion of mistakes for top- $k$  ranking. We also derive the optimal signal to noise ratio condition for the exact recovery of the top- $k$  set. The maximum likelihood estimator (MLE) is shown to achieve both optimal partial recovery and optimal exact recovery. On the other hand, we show another popular algorithm, the spectral method, is in general suboptimal. Our results complement the recent work (*Ann. Statist.* **47** (2019) 2204–2235) that shows both the MLE and the spectral method achieve the optimal sample complexity for exact recovery. It turns out the leading constants of the sample complexity are different for the two algorithms. Another contribution that may be of independent interest is the analysis of the MLE without any penalty or regularization for the BTL model. This closes an important gap between theory and practice in the literature of ranking.

**1. Introduction.** Given partially observed pairwise comparison data from  $n$  players, a central statistical question is how to optimally aggregate the comparison results and to find the leading top  $k$  players. This problem is known as top- $k$  ranking, which has important applications in many areas such as web search [11, 12] and competitive sports [19, 22]. In this paper, our goal is to study the statistical limits of both *partial* and *exact* recovery of the top- $k$  ranking problem.

We will focus on the popular Bradley–Terry–Luce (BTL) pairwise comparison model [3, 18]. That is, we observe  $L$  games played between  $i$  and  $j$ , and the outcome is modeled by

$$(1) \quad y_{ijl} \stackrel{\text{ind}}{\sim} \text{Bernoulli}\left(\frac{w_i^*}{w_i^* + w_j^*}\right), \quad l = 1, \dots, L.$$

We only observe outcomes from a small subset of pairs. This subset  $\mathcal{E}$  is modeled by edges generated by an Erdős–Rényi [13] random graph with connection probability  $p$  on the  $n$  players. More details of the model will be given in Section 2. With the observations  $\{y_{ijl}\}_{(i,j) \in \mathcal{E}, l \in [L]}$ , the goal is to reliably recover the set of top- $k$  players with the largest skill parameters  $w_i^*$ .

Theoretical properties of the top- $k$  ranking problem have been studied by [6–8, 14, 15, 21, 23] and references therein. The literature is mainly focused the problem of exact recovery. That is, to investigate the signal to noise ratio condition under which one can recovery the top- $k$  set without any error in probability. For this purpose, the state-of-the-art result is obtained by the recent work [7]. It was shown by [7] that both the MLE and the spectral method can

---

Received July 2021; revised December 2021.

MSC2020 subject classifications. 62F07.

Key words and phrases. Top- $k$  ranking, pairwise comparison, MLE, spectral method, partial recovery, exact recovery.

perfectly identify the top- $k$  players under optimal sample complexity up to some constant factor. This discovery was also verified by a numerical experiment that shows almost identical performances of the two methods. The results of [7] lead to the following intriguing research questions. What is the leading constant factor of the optimal sample complexity? Are the MLE and the spectral method still optimal if we take the leading constant into consideration?

In this paper, we give complete answers to the above questions. Our results show that while the MLE achieves a leading constant that is information-theoretically optimal, the spectral method only achieves a suboptimal constant. In particular, the MLE achieves exact recovery when

$$(2) \quad npL\Delta^2 > 2.001V(\kappa)(\sqrt{\log k} + \sqrt{\log(n-k)})^2,$$

and the spectral method requires

$$npL\Delta^2 > 2.001\bar{V}(\kappa)(\sqrt{\log k} + \sqrt{\log(n-k)})^2.$$

In the above two formulas,  $\Delta$  is the logarithmic gap of the skill parameters between the top- $k$  group and the rest of the players. The parameter  $\kappa$  is the dynamic range of the skill vector that will be defined in Section 2. The performances of the two methods are precisely characterized by the two functions  $V(\kappa)$  and  $\bar{V}(\kappa)$ , which are understood to be the effective variances of the two algorithms. The two functions satisfy the strict inequality that  $\bar{V}(\kappa) > V(\kappa)$  for all  $\kappa > 0$ , and the equality  $\bar{V}(\kappa) = V(\kappa)$  only holds when  $\kappa = 0$ . We also establish an information-theoretic lower bound that shows the MLE constant  $V(\kappa)$  is optimal, and it characterizes the phase transition boundary of exact recovery for the top- $k$  ranking problem.

We would like to emphasize that our results do not contradict the conclusions of [7]. On the contrary, the current paper complements and refines the results of [7]. The optimality claim made by [7] on both the MLE and the spectral method only refers to the order of the sample complexity. Our results show that the performances of the two algorithms can be drastically different when the dynamic range parameter  $\kappa$  is strictly positive. We are also able to explain why the numerical experiment conducted in [7] demonstrates nearly identical performances of the MLE and the spectral method. Note that the experiment in [7] was conducted with the skill parameters  $w_i^*$  only taking two possible values,  $e^\Delta$  or 1, depending on whether  $i$  belongs to the top- $k$  group or not. We show in Section 5 that this configuration of  $w^*$  is asymptotically equivalent to  $\kappa = 0$ , which is the only case that makes  $\bar{V}(\kappa) = V(\kappa)$ , and thus the nearly identical performances of the two algorithms are actually well expected by our theory. As long as  $w^*$  deviates from this simple two-piece structure, our extensive numerical experiments in this paper show that the MLE always performs better than the spectral method, and the advantage of the MLE is usually quite significant.

Another popular ranking algorithm is to select the  $k$  players with the most number of games won. This is often referred to as the Borda count [2] or the Copeland count [10]. This simple method does not require an explicit parametric model assumption, which makes it robust under potential model misspecification, an advantage over the MLE and the spectral method. The performance of the count method has been analyzed by [23] under a very general nonparametric setting. However, under the BTL model, it can be shown that the count method has an undesirable property that its error probability does not decrease to zero as  $L$  tends to infinity, which implies that it is not even rate-optimal under the BTL model. See Appendix F of the Supplementary Material [5] for an analysis of the count method.

In addition to the exact recovery results, we have also obtained a series of results for partial recovery. We observe that top- $k$  ranking can be viewed as a clustering problem. That is, one wants to cluster the players into two groups of sizes  $k$  and  $n - k$ , respectively. Therefore, it is more natural to consider the problem of partial recovery by analyzing the proportion of

players that are clustered into a wrong group. Clearly, this problem is more relevant in practice, since one rarely expects any real application where top- $k$  ranking can be done without any error. From a mathematical point of view, the partial recovery problem is more general and we will show in Section 3 that an optimal partial recovery error bound will lead to the optimal exact recovery condition (2). To the best of our knowledge, a systematic study of partial recovery for top- $k$  ranking has never been done in the literature. Our paper is perhaps the first work that formulates the top- $k$  ranking problem into a decision-theoretic framework and derives the minimax optimal partial recovery error rate. Similar to the results of exact recovery, we show that the MLE is also optimal for partial recovery. It has an exponential error bound with respect to a normalized Hamming loss. The error exponent is shown to depend on the variance function  $V(\kappa)$ . In comparison, the spectral method still achieves a suboptimal error rate for partial recovery, with the error exponent depending on  $\overline{V}(\kappa)$ .

Recently, a few papers provide sharp analysis of spectral methods on some high-dimensional estimation problems and show spectral methods can achieve optimal theoretical guarantees just as MLEs. For example, it was shown by [1] that spectral clustering achieves optimal community detection for a special class of stochastic block models (SBMs). The paper [16] proved spectral clustering is also optimal under Gaussian mixture models. We emphasize that the results of both papers imply that not only the order of the sample complexity of spectral clustering is optimal, but even the leading constant is optimal, at least in the setting of SBMs and Gaussian mixture models. The results of the current paper, however, show that the optimality of spectral methods may not hold under more complicated settings such as the BTL model.

Finally, we discuss another contribution of the paper that may be of independent interest. That is, we are able to give a sharp analysis of the MLE under the BTL model. Previous analyses of the MLE in the literature [7, 8, 21] all impose some additional regularization to address the challenge that the Hessian of the log-likelihood function is not well behaved. Whether the *vanilla* MLE works theoretically without any penalty or regularization remains an open problem. Our analysis solves this open problem by relating a regularized MLE to an  $\ell_\infty$ -constrained MLE. This allows us to show that the solution to the  $\ell_\infty$ -constrained MLE lies in the interior of the constraint. Thus, we can conclude that the  $\ell_\infty$ -constrained MLE is equivalent to the vanilla MLE in its original form. This equivalence then leads to the desired control of the spectrum of the Hessian matrix, which is the most critical step of our analysis.

The rest of the paper is organized as follows. We introduce the setting of the problem in Section 2. The results of the MLE and the spectral method will be given in Section 3 and Section 4, respectively. We then comprehensively compare the two methods in Section 5 by numerical experiments. Section 6 presents a minimax lower bound for partial recovery. In Section 7, we analyze the error rates of the MLE and the spectral method for each individual parameter. Due to the limit on pages, we include the proofs of the results of the MLE in Section 8 and the proofs of all the other results in the Supplementary Material.

We close this section by introducing some notation that will be used in the paper. For an integer  $d$ , we use  $[d]$  to denote the set  $\{1, 2, \dots, d\}$ . Given two numbers  $a, b \in \mathbb{R}$ , we use  $a \vee b = \max(a, b)$  and  $a \wedge b = \min(a, b)$ . We also write  $a_+ = \max(a, 0)$ . For two positive sequences  $\{a_n\}, \{b_n\}$ ,  $a_n \lesssim b_n$  or  $a_n = O(b_n)$  means  $a_n \leq Cb_n$  for some constant  $C > 0$  independent of  $n$ ,  $a_n = \Omega(b_n)$  means  $b_n = O(a_n)$ , and  $a_n \asymp b_n$  means  $a_n \lesssim b_n$  and  $b_n \lesssim a_n$ . We also write  $a_n = o(b_n)$  when  $\limsup_n \frac{a_n}{b_n} = 0$ . For a set  $S$ , we use  $\mathbb{I}\{S\}$  to denote its indicator function and  $|S|$  to denote its cardinality. For a vector  $v \in \mathbb{R}^d$ , its norms are defined by  $\|v\|_1 = \sum_{i=1}^d |v_i|$ ,  $\|v\|^2 = \sum_{i=1}^d v_i^2$  and  $\|v\|_\infty = \max_{1 \leq i \leq d} |v_i|$ . The notation  $\mathbb{1}_d$  means a  $d$ -dimensional column vector of all ones. For any  $v \in \mathbb{R}^d$ , we write  $\text{ave}(v) = d^{-1} \mathbb{1}_d^T v$ . Given  $p, q \in (0, 1)$ , the Kullback–Leibler divergence is defined by  $D(p\|q) = p \log \frac{p}{q} + (1 - p) \log \frac{1-p}{1-q}$ . For a natural number  $n$ ,  $\mathfrak{S}_n$  is the set of permutations on  $[n]$ . The notation  $\mathbb{P}$  and

$\mathbb{E}$  are used for generic probability and expectation whose distribution is determined from the context.

## 2. Models and methods.

**2.1. The BTL model.** We start by introducing the setting of our problem. Consider  $n$  players, and each one is associated with a positive latent skill parameter  $w_i^*$  for  $i \in [n]$ . The comparison scheme of the  $n$  players is characterized by an Erdős–Rényi random graph  $A \sim \mathcal{G}(n, p)$ . That is,  $A_{ij} \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$  for all  $1 \leq i < j \leq n$ . For a pair  $(i, j)$  that is connected by the random graph and  $A_{ij} = 1$ , we observe  $L$  games played between  $i$  and  $j$ . The outcome of the games is modeled by the Bradley–Terry–Luce (BTL) model (1). Our goal is to identify the top- $k$  players whose skill parameters  $w_i^*$ ’s have the largest values.

To formulate this problem from a decision-theoretic point of view, we reparametrize the BTL model (1) by a sorted vector  $\theta^*$  and a rank vector  $r^*$ . A sorted vector  $\theta^*$  satisfies  $\theta_1^* \geq \theta_2^* \geq \dots \geq \theta_n^*$ , and a rank vector  $r^*$  is an element of permutation  $r^* \in \mathfrak{S}_n$ . Then the BTL model (1) can be equivalently written as

$$(3) \quad y_{ijl} \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\psi(\theta_{r_i^*}^* - \theta_{r_j^*}^*)), \quad l = 1, \dots, L,$$

where  $\psi(\cdot)$  is the sigmoid function  $\psi(t) = \frac{1}{1+e^{-t}}$ . In the original representation, we have  $w_i^* = \exp(\theta_{r_i^*}^*)$  for all  $i \in [n]$ . With (3), the top- $k$  ranking problem is to identify the subset  $\{i \in [n] : r_i^* \leq k\}$  from the random comparison data. This is a typical semiparametric problem because of the presence of the nuisance parameter  $\theta^*$ .

**2.2. Loss function for top- $k$  ranking.** Our goal is to study optimal top- $k$  ranking in terms of both *partial* and *exact* recovery. We thus introduce a loss function to quantify the error of top- $k$  ranking. Given any  $\hat{r}, r^* \in \mathfrak{S}_k$ , define the normalized Hamming distance by

$$(4) \quad H_k(\hat{r}, r^*) = \frac{1}{2k} \left( \sum_{i=1}^n \mathbb{I}\{\hat{r}_i > k, r_i^* \leq k\} + \sum_{i=1}^n \mathbb{I}\{\hat{r}_i \leq k, r_i^* > k\} \right).$$

The definition (4) gives a natural loss function for top- $k$  ranking, since  $H_k(\hat{r}, r^*)$  can be equivalently written as the cardinality of the symmetric difference of the sets  $\{i \in [n] : \hat{r}_i \leq k\}$  and  $\{i \in [n] : r_i^* \leq k\}$  normalized by  $2k$ . The value of  $H_k(\hat{r}, r^*)$  is always within the unit interval  $[0, 1]$ . Moreover,  $H_k(\hat{r}, r^*) = 0$  if and only if  $\{i \in [n] : \hat{r}_i \leq k\} = \{i \in [n] : r_i^* \leq k\}$ .

The loss function (4) can be related to various quantities previously defined in the literature. One of the most popular distances to compare two rank vectors is the *Kendall tau distance* defined as

$$K(\hat{r}, r^*) = \frac{1}{n} \sum_{1 \leq i < j \leq n} \mathbb{I}\{\text{sign}(\hat{r}_i - \hat{r}_j) \text{sign}(r_i^* - r_j^*) < 0\}.$$

Since  $K(\hat{r}, r^*)$  counts all pairwise differences in the ranking relation, it is a stronger distance than (4). While  $K(\hat{r}, r^*) = 0$  requires  $\hat{r} = r^*$ ,  $H_k(\hat{r}, r^*) = 0$  only requires the two top- $k$  sets are identical regardless of the actual ranks of the members of the sets. In fact, the study of the BTL model under  $K(\hat{r}, r^*)$ , called full ranking, is also a very interesting problem, and will be considered in a different paper.

As we have discussed in Section 1, the top- $k$  ranking problem can be thought of as a special variable selection problem. Variable selection under the normalized Hamming loss has recently been studied by [4, 20]. Consider either a Gaussian sequence model or a regression

model with coefficient vector  $\beta^* \in \mathbb{R}^p$  that satisfies either  $\beta_j^* = 0$  or  $|\beta_j^*| > a$ . The papers [4, 20] consider estimating  $\beta^*$  under the loss

$$\bar{H}_s(\hat{\beta}, \beta^*) = \frac{1}{2s} \left( \sum_{j=1}^p \mathbb{I}\{|\hat{\beta}_j| > a, \beta_j^* = 0\} + \sum_{j=1}^p \mathbb{I}\{\hat{\beta}_j = 0, |\beta_j^*| > a\} \right),$$

where  $s$  is the number of  $\beta_j^*$ 's that are not zero. One can clearly see the similarity between the two loss functions  $H_k(\hat{r}, r^*)$  and  $\bar{H}_s(\hat{\beta}, \beta^*)$ . Similarly, the loss  $\bar{H}_s(\hat{\beta}, \beta^*)$  only characterizes the estimation error of the set  $\{j \in [p] : |\beta_j^*| > a\}$ , and  $\bar{H}_s(\hat{\beta}, \beta^*) = 0$  if and only if  $\{j \in [p] : |\hat{\beta}_j| > a\} = \{j \in [p] : |\beta_j^*| > a\}$ .

**2.3. Parameter space.** For the nuisance parameter  $\theta^*$  of the model (3), it is necessary that there exists a positive gap between  $\theta_k^*$  and  $\theta_{k+1}^*$  for the top- $k$  set  $\{i \in [n] : r_i^* \leq k\}$  to be identifiable. We introduce a parameter space for this purpose. For any  $0 \leq \Delta \leq \kappa$ , define

$$\Theta(k, \Delta, \kappa) = \{\theta \in \mathbb{R}^n : \theta_1 \geq \dots \geq \theta_n, \theta_k - \theta_{k+1} \geq \Delta, \theta_1 - \theta_n \leq \kappa\}.$$

For any  $\theta^* \in \Theta(k, \Delta, \kappa)$ , a positive  $\Delta$  guarantees that there is a separation between the group of top- $k$  players and the rest. The number  $\kappa$  is called *dynamic range* of the problem.<sup>1</sup> This is a very important quantity, since it is closely related to the effective variance of the problem. Our results will give the exact dependence of the top- $k$  ranking error on both  $\Delta$  and  $\kappa$ .

**2.4. MLE and spectral method.** We study and compare the performances of two algorithms in the paper. The first algorithm is based on the maximum likelihood estimator (MLE). For each  $(i, j)$ , we use the notation  $\bar{y}_{ij} = \frac{1}{L} \sum_{l=1}^L y_{ijl}$ . Throughout the paper, we adopt the convention of notation that  $A_{ij} = A_{ji}$  and  $\bar{y}_{ij} = 1 - \bar{y}_{ji}$ . Then the negative log-likelihood function is given by

$$(5) \quad \ell_n(\theta) = \sum_{1 \leq i < j \leq n} A_{ij} \left[ \bar{y}_{ij} \log \frac{1}{\psi(\theta_i - \theta_j)} + (1 - \bar{y}_{ij}) \log \frac{1}{1 - \psi(\theta_i - \theta_j)} \right].$$

Define the MLE,

$$(6) \quad \hat{\theta} \in \underset{\theta: \mathbb{1}_n^T \theta = 0}{\operatorname{argmin}} \ell_n(\theta).$$

It can be shown that  $\hat{\theta}$  is unique as long as the comparison graph is connected. Then set  $\hat{r}$  to be the rank of players based on  $\hat{\theta}$ . In other words, find any  $\hat{r} \in \mathfrak{S}_n$  such that  $\hat{\theta}_{\hat{\sigma}_1} \geq \dots \geq \hat{\theta}_{\hat{\sigma}_n}$  is satisfied, where  $\hat{\sigma}$  is the inverse of  $\hat{r}$ . We emphasize that the MLE (6) is written in its vanilla version, without any constraint or penalty. To the best of our knowledge, (6) has not been previously analyzed in the literature.

Another popular algorithm for ranking is the spectral method, also known as rank centrality proposed by [21]. Define a matrix  $P \in \mathbb{R}^{n \times n}$  by

$$(7) \quad P_{ij} = \begin{cases} \frac{1}{d} A_{ij} \bar{y}_{ji}, & i \neq j, \\ 1 - \frac{1}{d} \sum_{l \in [n] \setminus \{i\}} A_{il} \bar{y}_{li}, & i = j, \end{cases}$$

where  $d$  needs to be at least the maximum degree of the random graph  $A$ . We just set  $d = 2np$  throughout the paper. One can check that  $P$  is a transition matrix of a Markov chain. To see

<sup>1</sup>For readers who are familiar with [7], we note that our definitions of  $\Delta$  and  $\kappa$  are slightly different from those in [7].

why  $P$  is useful, we can compute the conditional expectation of  $P$  given the random graph  $A$ ,

$$P_{ij}^* = \begin{cases} \frac{1}{d} A_{ij} \psi(\theta_{r_j}^* - \theta_{r_i}^*), & i \neq j, \\ 1 - \frac{1}{d} \sum_{l \in [n] \setminus \{i\}} A_{il} \psi(\theta_{r_l}^* - \theta_{r_i}^*), & i = j. \end{cases}$$

The stationary distribution induced by the Markov chain  $P^*$  is

$$(\pi^*)^T = \left( \frac{\exp(\theta_{r_1}^*)}{\sum_{i=1}^n \exp(\theta_{r_i}^*)}, \dots, \frac{\exp(\theta_{r_n}^*)}{\sum_{i=1}^n \exp(\theta_{r_i}^*)} \right).$$

One can easily check that  $(\pi^*)^T P^* = (\pi^*)^T$ . Since  $\pi^*$  preserves the order of  $\{\theta_{r_i}^*\}$ , the set with the  $k$  largest  $\pi_i^*$ 's is the top- $k$  group. With the sample version  $P$ , we can first compute its stationary distribution  $\hat{\pi}$ , and then find any  $\hat{r} \in \mathfrak{S}_n$  such that  $\hat{\pi}_{\hat{\sigma}_1} \geq \dots \geq \hat{\pi}_{\hat{\sigma}_n}$ , with  $\hat{\sigma}$  being the inverse of  $\hat{r}$ .

**3. Results for the MLE.** We study the property of MLE in this section. Our first result gives theoretical guarantees for (6) under both  $\ell_2$  and  $\ell_\infty$  loss functions.

**THEOREM 3.1.** *Assume  $p \geq c_0 \frac{\log n}{n}$  for some sufficiently large constant  $c_0 > 0$  and  $\kappa \leq c_1$  for some constant  $c_1 > 0$ . Then, for the estimator  $\hat{\theta}$  defined by (6), we have*

$$(8) \quad \sum_{i=1}^n (\hat{\theta}_i - \theta_{r_i}^*)^2 \leq C \frac{1}{pL},$$

$$(9) \quad \max_{i \in [n]} |\hat{\theta}_i - \theta_{r_i}^*| \leq C \frac{\log n}{npL},$$

for some constant  $C > 0$  only depending on  $c_1$  with probability at least  $1 - O(n^{-7})$  uniformly over all  $r^* \in \mathfrak{S}_n$  and all  $\theta^* \in \Theta(k, 0, \kappa)$  such that  $\mathbf{1}_n^T \theta^* = 0$ .

Let us give some comments on the assumptions and conclusions of Theorem 3.1. We have established that the MLE achieves the error rates  $O(\frac{1}{pL})$  and  $O(\frac{\log n}{npL})$  for the squared  $\ell_2$  loss and the squared  $\ell_\infty$  loss, respectively. Both error rates are known to be optimal in the literature [7, 21]. Since the BTL model (3) is defined through pairwise differences of  $\theta_i^*$ 's, the model parameter is only identifiable up to a constant shift. We therefore require both  $\mathbf{1}_n^T \hat{\theta} = 0$  and  $\mathbf{1}_n^T \theta^* = 0$  so that the two vectors are properly aligned. Note that the results for parameter estimation do not need a positive  $\Delta$ , and we only assume  $\theta^* \in \Theta(k, 0, \kappa)$ . The condition  $p \geq c_0 \frac{\log n}{n}$  is imposed for the random graph  $A$  to be well behaved in terms of both its degrees and the eigenvalues of the graph Laplacian. In fact,  $p \gtrsim \frac{\log n}{n}$  is necessary to ensure the random graph is connected. Otherwise, ranking and parameter estimation would be impossible due to the identifiability issue caused by the lack of comparison between disconnected graph components. In the rest of the paper, some of the results will require a slightly stronger condition  $\frac{np}{\log n} \rightarrow \infty$ , but we will give very detailed remarks on when and why it will be needed. Last but not least, we require that the dynamic range  $\kappa$  to be bounded by a constant. One can certainly allow  $\kappa$  to tend to infinity, but the rates (8) and (9) would depend on  $\kappa$  exponentially [7, 21]. This is because the eigenvalues of the Hessian of the objective function of (6) will be exponentially small when  $\kappa$  diverges. In fact, when  $\kappa \rightarrow \infty$ , it is not clear whether MLE still leads to optimal error rates for parameter estimation. In this paper, we will

focus on the case  $\kappa = O(1)$ . We will see in later theorems that even with  $\kappa = O(1)$ , the exact value of  $\kappa$  still plays a fundamental role in top- $k$  ranking.

To the best of our knowledge, Theorem 3.1 is the first result in the literature that gives optimal rates for parameter estimation by *vanilla* MLE under the BTL model. Previous results in the literature including [7, 8, 21] all work with regularized MLE,

$$(10) \quad \hat{\theta}_\lambda = \operatorname{argmin}_{\theta: \mathbf{1}_n^T \theta = 0} \left[ \ell_n(\theta) + \frac{\lambda}{2} \|\theta\|^2 \right].$$

In particular, the recent paper [7] shows that  $\hat{\theta}_\lambda$  also achieves the optimal rates (8) and (9) for a  $\lambda$  that is chosen appropriately, though in practice it is known that the vanilla MLE performs very well. Theorem 3.1 shows that penalty is not needed for the MLE to be optimal, thus closing a gap between theory and practice.

The proof of Theorem 3.1 is built upon the elegant leave-one-out technique in [7]. We first show that with a sufficiently small  $\lambda$ , a (suboptimal)  $\ell_\infty$  bound for  $\hat{\theta}_\lambda$  can be transferred to  $\hat{\theta}$ . Then we apply a leave-one-out argument to derive the optimal rates (8) and (9). We also note that our leave-one-out argument is actually different from the form used in [7]. While the leave-one-out argument in [7] is applied together with a gradient descent analysis, we do not need to follow this gradient descent analysis because of the  $\ell_\infty$  bound that has already been obtained. As a result, we are able to remove the additional technical assumption  $\log L = O(\log n)$  that is imposed in [7]. A detailed analysis of the MLE will be given in Section 8.

Next, we study the theoretical property of  $\hat{r}$ , the rank induced by the MLE  $\hat{\theta}$ . Without loss of generality, let us assume  $k \leq \frac{n}{2}$  throughout the paper. The case  $k > \frac{n}{2}$  can be dealt with by a symmetric bottom- $k$  ranking problem. Before presenting the error bound for the loss function  $H_k(\hat{r}, r^*)$ , we need to introduce a few notation. We first define the effective variance of the MLE by

$$(11) \quad V(\kappa) = \max_{\substack{\kappa_1 + \kappa_2 \leq \kappa \\ \kappa_1, \kappa_2 \geq 0}} \frac{n}{k\psi'(\kappa_1) + (n-k)\psi'(\kappa_2)}.$$

Recall that  $\psi(t) = \frac{1}{1+e^{-t}}$  is the sigmoid function so that  $\psi'(t) = \psi(t)\psi(-t)$ . Since  $\kappa = O(1)$ , we have  $V(\kappa) \asymp 1$ . Then the signal to noise ratio is defined by

$$\text{SNR} = \frac{npL\Delta^2}{V(\kappa)}.$$

Note that SNR is a function of  $n, k, p, L, \Delta$ , but we suppress the dependence for simplicity of notation. The following theorem shows that  $H_k(\hat{r}, r^*)$  has an exponential rate with SNR appearing in the exponent.

**THEOREM 3.2.** *Assume  $\frac{np}{\log n} \rightarrow \infty$  and  $\kappa \leq c_1$  for some constant  $c_1 > 0$ . Then, for the rank vector  $\hat{r}$  that is induced by the MLE (6), there exists some  $\delta = o(1)$ , such that*

$$(12) \quad H_k(\hat{r}, r^*) \leq C \exp \left( -\frac{1}{2} \left( \frac{\sqrt{(1-\delta)\text{SNR}}}{2} - \frac{1}{\sqrt{(1-\delta)\text{SNR}}} \log \frac{n-k}{k} \right)_+^2 \right),$$

for some constant  $C > 0$  only depending on  $c_1$  with probability  $1 - o(1)$  uniformly over all  $r^* \in \mathfrak{S}_n$  and all  $\theta^* \in \Theta(k, \Delta, \kappa)$ .

The error exponent of (12) is complicated. We present a special case of the bound when  $k \asymp n$  to help understand the result.

COROLLARY 3.1. Assume  $\frac{np}{\log n} \rightarrow \infty$ ,  $\kappa = O(1)$  and  $k \asymp n$ . Then, as long as  $\text{SNR} \rightarrow \infty$ , the rank vector  $\hat{r}$  induced by the MLE (6) satisfies

$$(13) \quad H_k(\hat{r}, r^*) \leq \exp\left(-\left(1 - o(1)\right)\frac{\text{SNR}}{8}\right),$$

with probability  $1 - o(1)$  uniformly over all  $r^* \in \mathfrak{S}_n$  and all  $\theta^* \in \Theta(k, \Delta, \kappa)$ .

Under the additional assumption  $k \asymp n$ , the top- $k$  ranking problem can be viewed as a clustering or community detection problem, with the goal to divide the  $n$  players into two groups of sizes  $k$  and  $n - k$ , respectively. The exponential convergence rate (13) is in a typical form of optimal clustering error [17, 25]. It is intuitively clear that a larger SNR leads to a faster convergence rate. When the sizes of the two clusters are of different orders, one can obtain a more general convergence rate in the form of (12). The extra term  $\log \frac{n-k}{k}$  characterizes the unbalancedness of the two clusters. We note that for variable selection under Hamming loss [4, 20], the optimal rate is very similar to the form of (12). This is because variable selection can also be thought of as clustering with two clusters of sizes  $s$  and  $p - s$ , whose orders can potentially be different.

Theorem 3.2 and Corollary 3.1 together reveal an interesting phenomenon for top- $k$  ranking. The result shows that the top- $k$  ranking problem can be very different for different orders of  $k$ . We note that in order to successfully identify the majority of the set  $\{i \in [n] : r_i^* \leq k\}$ , we need to have  $H_k(\hat{r}, r^*) \rightarrow 0$ . When  $k = n/4$ , Corollary 3.1 shows that  $H_k(\hat{r}, r^*) \rightarrow 0$  is achieved when  $\text{SNR} \rightarrow \infty$ . In comparison, when  $k = 5$ , Theorem 3.2 shows that  $H_k(\hat{r}, r^*) \rightarrow 0$  when  $\text{SNR} > (1 + \epsilon)2 \log n$  for some arbitrarily small constant  $\epsilon > 0$ . In other words, in terms of partial recovery consistency, top-quarter ranking is an easier problem than top-5 ranking. In general, a larger SNR is required for a smaller  $k$  according to the formula (12).

Compared with Theorem 3.1, we need a slightly stronger condition  $\frac{np}{\log n} \rightarrow \infty$  for Theorem 3.2 and Corollary 3.1. If we only assume  $p \geq c_0 \frac{\log n}{n}$ , the  $1 - \delta$  factor in the exponent of (12) can be replaced by  $1 - \epsilon$  with some  $\epsilon$  of constant order. The constant  $\epsilon$  can be made arbitrarily small as long as  $c_0$  is sufficiently large.

The proof of Theorem 3.2 relies on a very interesting lemma that is stated below.

LEMMA 3.1. Suppose  $\hat{r}$  is a rank vector induced by  $\hat{\theta}$ , we then have

$$H_k(\hat{r}, r^*) \leq \frac{1}{k} \min \left[ \sum_{i: r_i^* \leq k} \mathbb{I}\{\hat{\theta}_i \leq t\} + \sum_{i: r_i^* > k} \mathbb{I}\{\hat{\theta}_i \geq t\} \right].$$

The inequality holds for any  $r^* \in \mathfrak{S}_n$ .

We will prove Lemma 3.1 in Appendix E. This inequality shows that the error of ranking  $\hat{\theta}$  is bounded by the error of any thresholding rule. Using this result, we immediately obtain that

$$\mathbb{E} H_k(\hat{r}, r^*) \leq \frac{1}{k} \min \left[ \sum_{i: r_i^* \leq k} \mathbb{P}(\hat{\theta}_i \leq t) + \sum_{i: r_i^* > k} \mathbb{P}(\hat{\theta}_i \geq t) \right].$$

We then obtain the exponential error bound (12) by carefully analyzing the probability  $\mathbb{P}(\hat{\theta}_i \leq t)$  (or  $\mathbb{P}(\hat{\theta}_i \geq t)$ ) for each  $i \in [n]$ . The analysis of  $\mathbb{P}(\hat{\theta}_i \leq t)$  is quite involved. We need to first obtain a local linear expansion of the MLE at each coordinate, and then apply the leave-one-out technique introduced by [7] to decouple the dependence between the data and the coefficients of the local linear expansion. The details will be given in Section 8.

The result of Theorem 3.2 immediately implies a condition for exact recovery of the top- $k$  set. By the definition of  $H_k(\hat{r}, r^*)$ , it is easy to see that

$$(14) \quad H_k(\hat{r}, r^*) \in \{0, (2k)^{-1}, 2(2k)^{-1}, 3(2k)^{-1}, \dots, 1\}.$$

Then as long as  $H_k(\hat{r}, r^*) < (2k)^{-1}$ , we must have  $H_k(\hat{r}, r^*) = 0$ . Under the condition that the right-hand side of (12) is smaller than  $(2k)^{-1}$ , we obtain exact recovery of the top- $k$  set. This result is stated as follows.

THEOREM 3.3. Assume  $\frac{np}{\log n} \rightarrow \infty$ ,  $\kappa = O(1)$ , and

$$(15) \quad \frac{npL\Delta^2}{V(\kappa)} > (1 + \epsilon)2(\sqrt{\log k} + \sqrt{\log(n - k)})^2,$$

for some arbitrarily small constant  $\epsilon > 0$ . Then, for the rank vector  $\hat{r}$  that is induced by the MLE (6), we have  $H_k(\hat{r}, r^*) = 0$  with probability  $1 - o(1)$  uniformly over all  $r^* \in \mathfrak{S}_n$  and all  $\theta^* \in \Theta(k, \Delta, \kappa)$ .

We remark that the condition  $\frac{np}{\log n} \rightarrow \infty$  can be relaxed to  $p \geq c_0 \frac{\log n}{n}$  for a sufficiently large constant  $c_0$  without affecting the conclusion of Theorem 3.3. The result of Theorem 3.3 improves the exact recovery threshold obtained in the literature. The paper [7] proves that the MLE exactly recovers the top- $k$  set when  $npL\Delta^2 > C(\sqrt{\log k} + \sqrt{\log(n - k)})^2$  for some sufficiently large constant  $C > 0$ . We complement the result of [7] by showing that the leading constant should be  $2V(\kappa)$ , an increasing function of the dynamic range  $\kappa$ . Moreover, the symmetry of  $k$  and  $n - k$  in (15) agrees with the understanding that top- $k$  ranking and bottom- $k$  ranking are mathematically equivalent.

The next theorem shows that the exact recovery threshold (15) is optimal, and cannot be further improved.

THEOREM 3.4. Assume  $\frac{np}{\log n} \rightarrow \infty$ ,  $\kappa = O(1)$ ,  $(\log n)^8 = O(L)$  and

$$(16) \quad \frac{npL\Delta^2}{V(\kappa)} < (1 - \epsilon)2(\sqrt{\log k} + \sqrt{\log(n - k)})^2,$$

for some arbitrarily small constant  $\epsilon > 0$ . Then we have

$$\liminf_{n \rightarrow \infty} \inf_{\hat{r}} \sup_{\substack{r^* \in \mathfrak{S}_n \\ \theta^* \in \Theta(k, \Delta, \kappa)}} \mathbb{P}_{(\theta^*, r^*)}(H_k(\hat{r}, r^*) > 0) \geq 0.95,$$

where we use the notation  $\mathbb{P}_{(\theta^*, r^*)}$  for the data generating process (3).

The proof of Theorem 3.4 relies on a precise lower bound characterization of the maximum of dependent binomial random variables. The extra assumption  $L \gtrsim (\log n)^8$  allows us to apply a high-dimensional central limit theorem [9] for this purpose. Without this additional technical condition, we are not aware of any probabilistic tool to deal with maximum of dependent binomial random variables.

Theorem 3.3 and Theorem 3.4 together nail down the phase transition boundary of exact recovery, which is  $\frac{npL\Delta^2}{V(\kappa)} = 2(\sqrt{\log k} + \sqrt{\log(n - k)})^2$ . Thus, the MLE is an optimal procedure that achieves this boundary. The lower bound result of Theorem 3.4 also suggests that the partial recovery error rate obtained in Theorem 3.2 cannot be improved, since otherwise one would obtain a better SNR condition for exact recovery in Theorem 3.3. A rigorous minimax lower bound for partial recovery will be given in Section 6.

**4. Results for the spectral method.** In this section, we study the theoretical property of the spectral method, also known as rank centrality [21]. Let  $\hat{\pi}$  be the stationary distribution of the Markov chain with transition probability (7). The estimation error of  $\hat{\pi}$  has already been investigated by [7, 21]. For both  $\ell_2$  and  $\ell_\infty$  loss functions, it has been shown by [7] that  $\hat{\pi}$  achieves the optimal rates (8) and (9) after an appropriate scaling. We therefore directly study the accuracy of the rank vector  $\hat{r}$  induced by  $\hat{\pi}$ . This is where we can see the difference between the MLE and the spectral method.

We first define the effective variance of the spectral method,

$$(17) \quad \bar{V}(\kappa) = \max_{\substack{\kappa_1 + \kappa_2 \leq \kappa \\ \kappa_1, \kappa_2 \geq 0}} \frac{k\psi'(\kappa_1)(1 + e^{\kappa_1})^2 + (n - k)\psi'(\kappa_2)(1 + e^{-\kappa_2})^2}{(k\psi(\kappa_1) + (n - k)\psi(-\kappa_2))^2/n}.$$

Note that  $\bar{V}(\kappa) \asymp 1$  when  $\kappa = O(1)$ . The signal to noise ratio is defined by

$$\overline{\text{SNR}} = \frac{npL\Delta^2}{\bar{V}(\kappa)}.$$

The error rate of the spectral method with respect to  $H_k(\hat{r}, r^*)$  is stated as follows.

**THEOREM 4.1.** *Assume  $\frac{np}{\log n} \rightarrow \infty$  and  $\kappa \leq c_1$  for some constant  $c_1 > 0$ . Then, for the rank vector  $\hat{r}$  that is induced by the stationary distribution of the Markov chain (7), there exists some  $\delta = o(1)$ , such that*

$$(18) \quad H_k(\hat{r}, r^*) \leq C \exp\left(-\frac{1}{2}\left(\frac{\sqrt{(1 - \delta)\overline{\text{SNR}}}}{2} - \frac{1}{\sqrt{(1 - \delta)\overline{\text{SNR}}}} \log \frac{n - k}{k}\right)_+^2\right),$$

for some constant  $C > 0$  only depending on  $c_1$  with probability  $1 - o(1)$  uniformly over all  $r^* \in \mathfrak{S}_n$  and all  $\theta^* \in \Theta(k, \Delta, \kappa)$ .

The formula (18) characterizes the convergence rate of partial recovery of the top- $k$  set by the spectral method. It can be compared with the MLE error bound (12). The only difference lies in the effective variance of the two methods. We will show in Lemma 5.1 that  $\bar{V}(\kappa) \geq V(\kappa)$  and the equality only holds when  $\kappa = 0$ . Therefore, the spectral method is not optimal in general. Detailed comparisons of the two algorithms will be given in Section 5.

By the property (14), we immediately obtain an exact recovery result from Theorem 4.1.

**THEOREM 4.2.** *Assume  $\frac{np}{\log n} \rightarrow \infty$ ,  $\kappa = O(1)$ , and*

$$(19) \quad \frac{npL\Delta^2}{\bar{V}(\kappa)} > (1 + \epsilon)2(\sqrt{\log k} + \sqrt{\log(n - k)})^2,$$

for some arbitrarily small constant  $\epsilon > 0$ . Then, for the rank vector  $\hat{r}$  that is induced by the stationary distribution of the Markov chain (7), we have  $H_k(\hat{r}, r^*) = 0$  with probability  $1 - o(1)$  uniformly over all  $r^* \in \mathfrak{S}_n$  and all  $\theta^* \in \Theta(k, \Delta, \kappa)$ .

It has been shown in [7] that the spectral method exactly recovers the top- $k$  set when  $npL\Delta^2 > C(\sqrt{\log k} + \sqrt{\log(n - k)})^2$  for some sufficiently large constant  $C > 0$ . Without specifying the constant  $C$ , one cannot tell the difference between the MLE and the spectral method. In view of the lower bound result given by Theorem 3.4, the exact recovery threshold (19) of the spectral method does not achieve the phase transition boundary for a general  $\kappa$ . A careful reader may wonder whether this is resulted from a loose analysis in the proof. Our next result shows that the suboptimality of the spectral method is intrinsic.

THEOREM 4.3. Assume  $\frac{np}{\log n} \rightarrow \infty$ ,  $\kappa \leq c_1$  for some constant  $c_1 > 0$ ,  $k \rightarrow \infty$  and

$$(20) \quad \frac{npL\Delta^2}{\bar{V}(\kappa)} < (1 - \epsilon)2(\sqrt{\log k} + \sqrt{\log(n - k)})^2,$$

for some arbitrarily small constant  $\epsilon > 0$ . Then, for the rank vector  $\hat{r}$  that is induced by the stationary distribution of the Markov chain (7), we have

$$\liminf_{n \rightarrow \infty} \sup_{\substack{r^* \in \mathfrak{S}_n \\ \theta^* \in \Theta(k, \Delta, \kappa)}} \mathbb{P}_{(\theta^*, r^*)}(\mathbf{H}_k(\hat{r}, r^*) > 0) \geq 0.95.$$

Moreover, there exists some  $\delta = o(1)$ , such that

$$(21) \quad \sup_{\substack{r^* \in \mathfrak{S}_n \\ \theta^* \in \Theta(k, \Delta, \kappa)}} \mathbb{E}_{(\theta^*, r^*)} \mathbf{H}_k(\hat{r}, r^*) \geq C \exp\left(-\frac{1}{2}\left(\frac{\sqrt{(1 + \delta)\text{SNR}}}{2} - \frac{1}{\sqrt{(1 + \delta)\text{SNR}}} \log \frac{n - k}{k}\right)_+^2\right),$$

for some constant  $C > 0$  only depending on  $c_1$  and  $\epsilon$ .

Theorem 4.3 shows that the results of Theorem 4.1 and Theorem 4.2 on the performance of spectral method are sharp, under the additional condition that  $k \rightarrow \infty$ . The conclusion of Theorem 4.3 can also be extended to the case of  $k = O(1)$  via a similar argument that is used in the proof of Theorem 3.4, as long as the technical condition  $(\log n)^8 = O(np)$  is further imposed.

To close this section, we remark that all the theorems we have obtained for the spectral method can be stated under the weaker assumption  $p \geq c_0 \frac{\log n}{n}$  for some sufficiently large constant  $c_0 > 0$ , as long as the  $\delta$  in (18) and (21) are replaced by some sufficiently small constant.

**5. Comparison of the two methods.** In this section, we compare the MLE and the spectral method based on the results obtained in Section 3 and Section 4. The statistical properties of the two methods in terms of partial and exact recovery are characterized by the two variance functions  $V(\kappa)$  and  $\bar{V}(\kappa)$ , respectively. We first give a direct comparison of the two functions by plotting them together with different values of  $k/n$ . We observe in Figure 1 that  $\bar{V}(\kappa) \geq V(\kappa)$  for all  $\kappa \geq 0$ . This inequality is rigorously established by the following lemma.

LEMMA 5.1. For  $V(\kappa)$  and  $\bar{V}(\kappa)$  defined in (11) and (17), respectively, we have

$$\bar{V}(\kappa) \geq V(\kappa),$$

for all  $\kappa \geq 0$ . Moreover, the equality holds if and only if  $\kappa = 0$ .

PROOF. By Jensen's inequality, we have

$$(22) \quad \frac{k \frac{e^{\kappa_1}}{(1+e^{\kappa_1})^2} + (n-k) \frac{e^{-\kappa_2}}{(1+e^{-\kappa_2})^2}}{ke^{\kappa_1} + (n-k)e^{-\kappa_2}} \geq \left( \frac{k \frac{e^{\kappa_1}}{1+e^{\kappa_1}} + (n-k) \frac{e^{-\kappa_2}}{1+e^{-\kappa_2}}}{ke^{\kappa_1} + (n-k)e^{-\kappa_2}} \right)^2.$$

Another way to see the above inequality is to construct a random variable  $X$  such that  $\mathbb{P}(X = \frac{1}{1+e^{\kappa_1}}) = \frac{ke^{\kappa_1}}{ke^{\kappa_1} + (n-k)e^{-\kappa_2}}$  and  $\mathbb{P}(X = \frac{1}{1+e^{-\kappa_2}}) = \frac{(n-k)e^{-\kappa_2}}{ke^{\kappa_1} + (n-k)e^{-\kappa_2}}$ . Then (22) is equivalent to  $\mathbb{E}X^2 \geq (\mathbb{E}X)^2$ . The inequality (22) can be rearranged into

$$(23) \quad \frac{k\psi'(\kappa_1)(1+e^{\kappa_1})^2 + (n-k)\psi'(\kappa_2)(1+e^{-\kappa_2})^2}{(k\psi(\kappa_1) + (n-k)\psi(-\kappa_2))^2/n} \geq \frac{n}{k\psi'(\kappa_1) + (n-k)\psi'(\kappa_2)}.$$

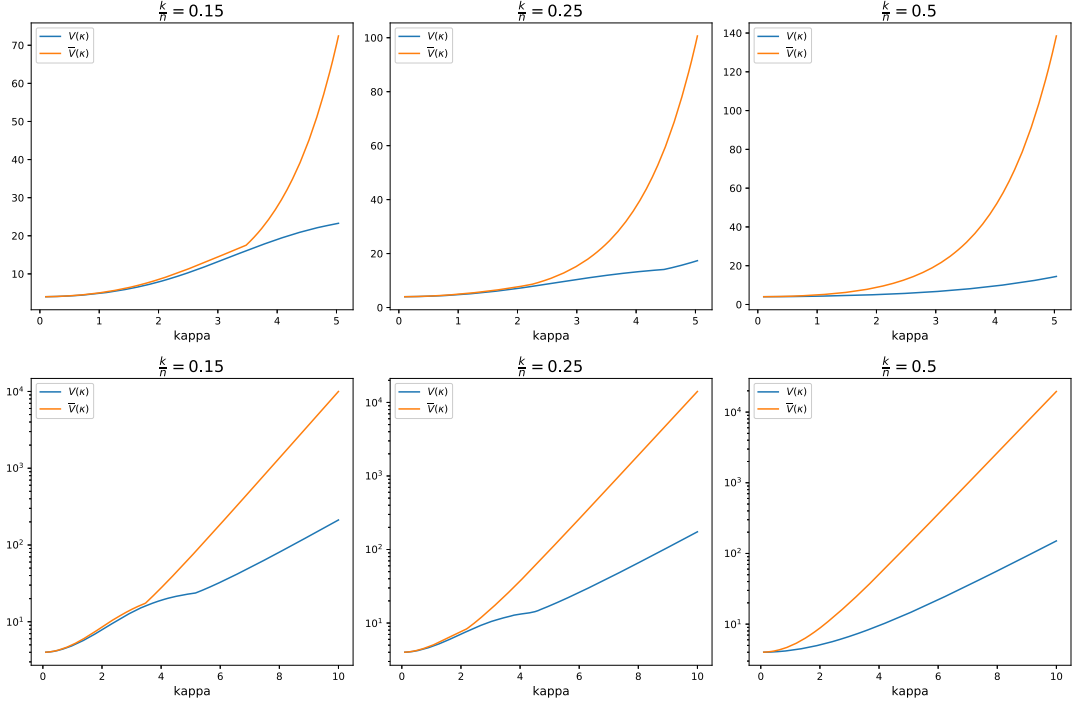


FIG. 1. The functions  $V(\kappa)$  and  $\bar{V}(\kappa)$  with  $k/n \in \{0.15, 0.25, 0.5\}$ . In the first row, we plot the functions for  $\kappa \in [0, 5]$ . The second row plots the same functions for  $\kappa \in [0, 10]$  in a logarithmic scale to better illustrate the global structure. It is very interesting that both  $V(\kappa)$  and  $\bar{V}(\kappa)$  have a point at which the derivative is not continuous. Before this critical point, the optimization of  $V(\kappa)$  is achieved by  $(\kappa_1^*, \kappa_2^*) = (0, \kappa)$ . Right after the critical point,  $\kappa_1^*$  is immediately bounded away from 0 and  $\kappa_2^*$  is immediately bounded away from  $\kappa$ . The same property also holds for  $\bar{V}(\kappa)$ . Moreover, the critical point occurs earlier as  $k/n$  becomes larger (when  $k/n \leq 1/2$ ).

Taking maximum over  $\kappa_1$  and  $\kappa_2$  on both sides, we obtain the inequality  $\bar{V}(\kappa) \geq V(\kappa)$ . When  $\kappa = 0$ , we obviously have  $V(\kappa) = \bar{V}(\kappa)$ . When  $\kappa > 0$ , we need to show  $V(\kappa) \neq \bar{V}(\kappa)$ . The optimization of  $V(\kappa)$  must be achieved by some  $(\kappa_1^*, \kappa_2^*) \neq (0, 0)$ . For such  $(\kappa_1^*, \kappa_2^*)$ , the constructed random variable  $X$  has a positive variance, and thus both inequalities (22) and (23) are strict. We then have

$$\begin{aligned} \bar{V}(\kappa) &\geq \frac{k\psi'(\kappa_1^*)(1 + e^{\kappa_1^*})^2 + (n - k)\psi'(\kappa_2^*)(1 + e^{-\kappa_2^*})^2}{(k\psi(\kappa_1^*) + (n - k)\psi(-\kappa_2^*))^2/n} \\ &> \frac{n}{k\psi'(\kappa_1^*) + (n - k)\psi'(\kappa_2^*)} \\ &= V(\kappa). \end{aligned}$$

The proof is complete.  $\square$

The comparison between  $V(\kappa)$  and  $\bar{V}(\kappa)$  shows that the spectral method is not optimal in general. It has a worse error exponent for partial recovery and requires a larger signal to noise ratio threshold for exact recovery. In fact, the difference  $\bar{V}(\kappa) - V(\kappa)$  eventually grows exponentially fast as a function of  $\kappa$ . See Figure 1.

Note that both  $V(\kappa)$  and  $\bar{V}(\kappa)$  are the worst-case effective variances with respect to the parameter space  $\Theta(k, \Delta, \kappa)$  for the two algorithms. In Section 7, we will further show that the MLE outperforms the spectral method for each  $\theta^* \in \Theta(k, \Delta, \kappa)$ . This conclusion is supported by extensive numerical experiments. We set  $n = 200$ ,  $p = 0.25$ ,  $L = 20$  and  $k = 50$  throughout the experiments.

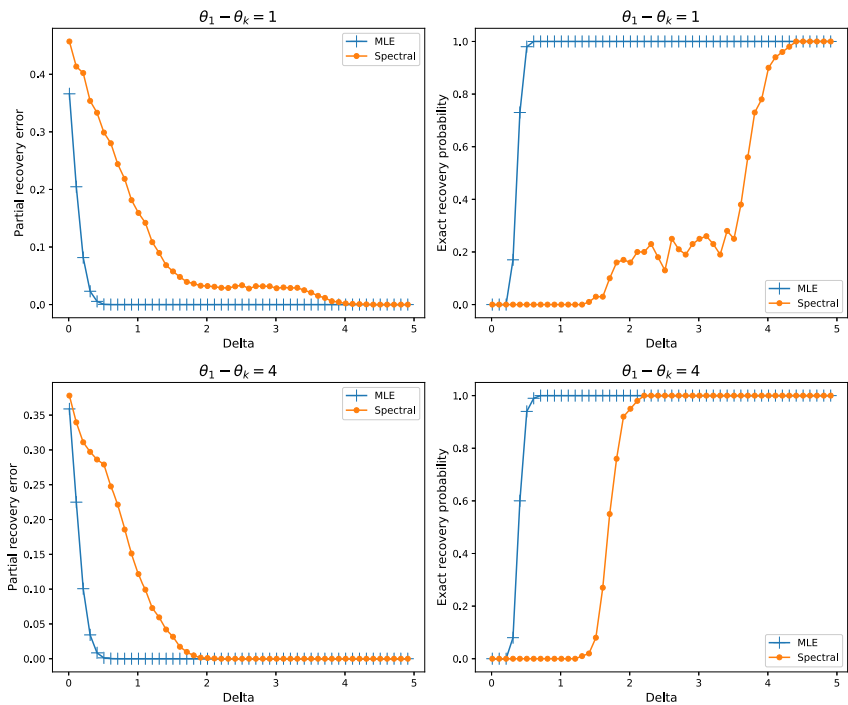


FIG. 2. The partial recovery error (left) and the exact recovery probability (right) for the MLE and the spectral method. The parameter  $\theta^*$  is chosen to be a piecewise constant vector of four pieces of sizes 25, 25, 75, 75. The plots are obtained by averaging 100 independent experiments.

In our first experiment, we consider  $\theta^* \in \mathbb{R}^n$  that has four pieces, with the three change-points located at  $\{25, 50, 200\}$ . The values of the four pieces are set as 10,  $10 - \tau$ ,  $10 - \tau - \Delta$  and 0, respectively, where  $\tau = \theta_1^* - \theta_k^* \in \{1, 4\}$  and  $\Delta$  is varied from 0.01 to 5. We apply both the MLE and the spectral method to the data. Figure 2 shows the results for both partial and exact recovery. We observe that the MLE consistently outperforms the spectral method.

In the second experiment, we consider  $\theta^* \in \mathbb{R}^n$  that has four pieces, with the three change-points located at  $\{50(1 - \rho), 50, 50 + 150\rho\}$ . The values of the four pieces are set as 10, 6,  $6 - \Delta$  and 0, respectively. The parameter  $\rho$  is chosen in  $\{0.1, 0.5, 0.9\}$  and  $\Delta$  is varied from 0.01 to 3. The performance of the two methods for partial and exact recovery are plotted in Figure 3. Again, the MLE always outperforms the spectral method.

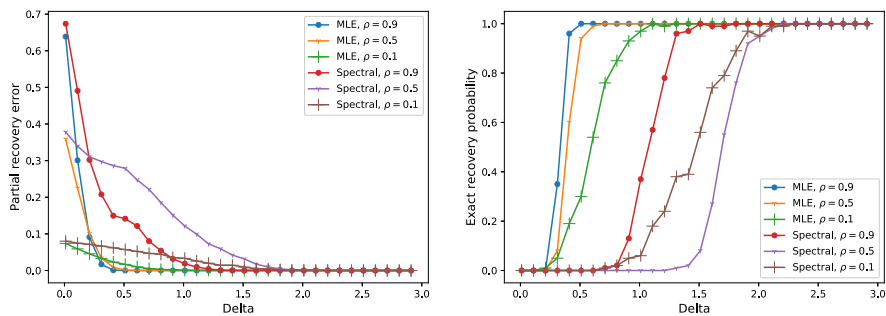


FIG. 3. The partial recovery error (left) and the exact recovery probability (right) for the MLE and the spectral method. The parameter  $\theta^*$  is chosen to be a piecewise constant vector of four pieces of sizes  $50(1 - \rho)$ ,  $50\rho$ ,  $150(1 - \rho)$ ,  $150\rho$ . The plots are obtained by averaging 100 independent experiments.

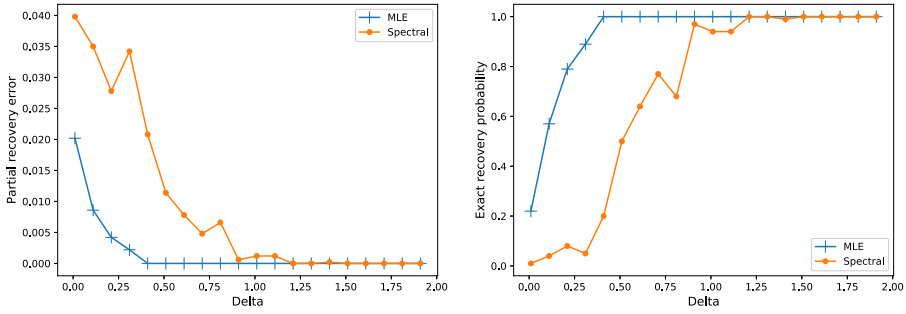


FIG. 4. The partial recovery error (left) and the exact recovery probability (right) for the MLE and the spectral method. The parameter  $\theta^*$  is randomly generated from some distribution. The plots are obtained by averaging 100 independent experiments.

Next, we consider a  $\theta^* \in \mathbb{R}^n$  that has a more complicated structure. We fix  $\theta_1^* = 10$ ,  $\theta_{200}^* = 0$ , generate  $\theta_2^*, \dots, \theta_{50}^*$  from Uniform[6, 10] and generate  $\theta_{51}^*, \dots, \theta_{199}^*$  from Uniform[0, 6 -  $\Delta$ ], and we vary  $\Delta$  from 0.01 to 2. We find even for such randomly generated  $\theta^*$ 's, the MLE always outperforms the spectral method. The results are summarized in Figure 4 for both partial and exact recovery.

In summary, we are able to confirm that the MLE is a much better algorithm than the spectral method under various scenarios. Our results complement the analysis in [7]. It is claimed in [7] that both the MLE and the spectral method are optimal in terms of the order of the exact recovery threshold. In addition, the paper conducts a very curious numerical experiment that shows the performances of the MLE and the spectral method are nearly identical. We note that the  $\theta^*$  chosen in the numerical experiment of [7] is a piecewise constant vector with only two pieces. We will explain why this choice leads to nearly identical performances of the two algorithms. Let us first conduct a similar experiment to replicate this conclusion. We continue to use the setting  $n = 200$ ,  $p = 0.25$ ,  $L = 20$  and  $k = 50$ . Then choose  $\theta^*$  such that  $\theta_1^* = \dots = \theta_{50}^* = \Delta$  and  $\theta_{51}^* = \dots = \theta_{200}^* = 0$ . Figure 5 plots the results of partial and exact recovery with  $\Delta$  varied from 0.01 to 0.55. For both partial recovery and exact recovery, the results are indeed nearly identical for the two algorithms. This phenomenon can be easily explained by our theory. For  $\theta^* \in \Theta(k, \Delta, \kappa)$  with only two pieces, we must have  $\kappa = \Delta$ . When  $\Delta = o(1)$ , we have  $V(\kappa) = (1 + o(1))V(0)$  and  $\bar{V}(\kappa) = (1 + o(1))\bar{V}(0)$ . This leads to the relation  $\bar{V}(\kappa) = (1 + o(1))V(\kappa)$ , and thus the spectral method has the same asymptotic error exponent for partial recovery and achieves the optimal phase transition boundary for exact recovery. When  $\Delta$  does not tend to zero but of a constant order, we have  $\overline{\text{SNR}} \gtrsim npL \gg \log n$ ,

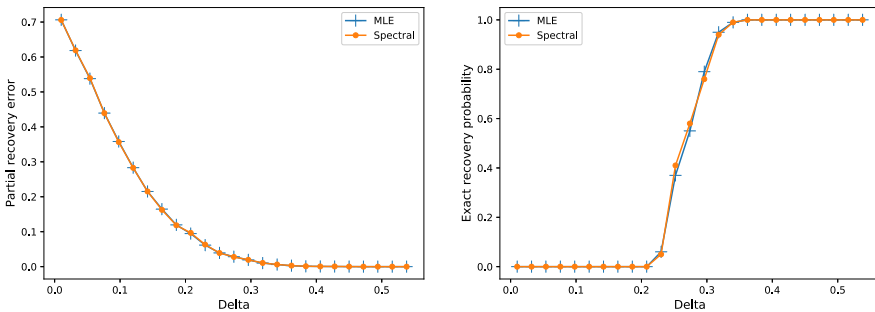


FIG. 5. The partial recovery error (left) and the exact recovery probability (right) for the MLE and the spectral method. The parameter  $\theta^*$  is chosen to be a piecewise constant vector of two pieces of sizes 50 and 150. The plots are obtained by averaging 100 independent experiments.

and the error bound (18) already leads to exact recovery because of the large value of  $\overline{\text{SNR}}$ . In either case, the spectral method is optimal. Let us summarize the optimality of the spectral method under this special situation by the following corollary.

**COROLLARY 5.1.** *Assume  $\frac{np}{\log n} \rightarrow \infty$  and  $\kappa = \Delta \leq c_1$  for some constant  $c_1 > 0$ . Then, for the rank vector  $\widehat{r}$  that is induced by the stationary distribution of the Markov chain (7), there exists some  $\delta = o(1)$ , such that*

$$H_k(\widehat{r}, r^*) \leq C \exp\left(-\frac{1}{2}\left(\frac{\sqrt{(1-\delta)\text{SNR}}}{2} - \frac{1}{\sqrt{(1-\delta)\text{SNR}}} \log \frac{n-k}{k}\right)_+^2\right),$$

for some constant  $C > 0$  only depending on  $c_1$  with probability  $1 - o(1)$  uniformly over all  $r^* \in \mathfrak{S}_n$  and all  $\theta^* \in \Theta(k, \Delta, \Delta)$ . Moreover, as long as

$$\frac{npL\Delta^2}{V(\kappa)} > (1 + \epsilon)2(\sqrt{\log k} + \sqrt{\log(n-k)})^2,$$

for some arbitrarily small constant  $\epsilon > 0$ . Then  $H_k(\widehat{r}, r^*) = 0$  with probability  $1 - o(1)$  uniformly over all  $r^* \in \mathfrak{S}_n$  and all  $\theta^* \in \Theta(k, \Delta, \Delta)$ .

To close this section, we remark that according to the equality condition of Lemma 5.1, the two-piece  $\theta^*$ , or equivalently  $\kappa = \Delta$ , is essentially the only situation where the spectral method is optimal and performs as well as the MLE. Moreover, since both functions  $V(\kappa)$  and  $\overline{V}(\kappa)$  are increasing, the setting with  $\kappa = \Delta$  leads to the smallest effective variance, and thus provides the two algorithms with the most favorable scenario.

**6. Minimax lower bound of partial recovery.** The purpose of this section is to show that the partial recovery error rate (12) achieved by the MLE cannot be improved from a minimax perspective. We are able to establish a matching lower bound for Theorem 3.2 using a slightly more general parameter space. Define

$$(24) \quad \Theta'(k, \Delta, \kappa) = \{\theta \in \mathbb{R}^n : \theta_1 \geq \cdots \geq \theta_n, \theta_k - \theta_{k+2} \geq \Delta, \theta_1 - \theta_n \leq \kappa\}.$$

Compared with  $\Theta(k, \Delta, \kappa)$ , the new definition (24) imposes a gap between  $\theta_k$  and  $\theta_{k+2}$ . It is clear that  $\Theta(k, \Delta, \kappa) \subset \Theta'(k, \Delta, \kappa)$ , and the only difference of  $\Theta'(k, \Delta, \kappa)$  is the ambiguity of  $\theta_{k+1}$ . The player ranked at the  $(k+1)$ th position does not necessarily has a gap from either the top group or the bottom group. Though this additional uncertainty clearly better models scenarios in many real applications of top- $k$  ranking, the main reason we adopt the slightly larger parameter space is to have a clean lower bound analysis. Directly establishing a lower bound for  $\Theta(k, \Delta, \kappa)$  is still possible, but it requires some additional technical assumptions that make the problem unnecessarily involved.

Throughout this section, we assume that (16) holds. This is the regime of partial recovery, since exact recovery is impossible by Theorem 3.4. We first remark that with a slight modification of the proof of Theorem 3.2, the MLE can be shown to achieve the same error rate (12) over the parameter space  $\Theta'(k, \Delta, \kappa)$  as well. Thus, the space  $\Theta'(k, \Delta, \kappa)$  does not increase the statistical complexity of the problem.

Our lower bound analysis is based on the two least favorable vectors  $\theta', \theta'' \in \Theta'(k, \Delta, \kappa)$ . They are constructed as follows. Let  $\rho = o(1)$  be a vanishing sequence that tends to zero with a sufficiently slow rate. We define  $\kappa_1^*$  and  $\kappa_2^*$  such that the optimization (11) is achieved at  $(\kappa_1, \kappa_2) = (\kappa_1^*, \kappa_2^*)$ . Then define  $\theta'_i = \kappa_1^*$  for  $1 \leq i \leq k - \rho k$ ,  $\theta'_i = 0$  for  $k - \rho k < i \leq k$ ,  $\theta'_i = -\Delta$  for  $k < i \leq k + \rho(n - k)$  and  $\theta'_i = -\kappa_2^*$  for  $k + \rho(n - k) < i \leq n$ . For  $\theta''$ , we let

$\theta''_i = \theta'_i$  for all  $i \in [n] \setminus \{k+1\}$  and set  $\theta''_{k+1} = 0$ . We will show that there exist  $r', r'' \in \mathfrak{S}_n$ , so that the hardness of top- $k$  ranking is characterized by an optimal testing problem,

$$(25) \quad \inf_{0 \leq \phi \leq 1} \left[ \mathbb{E}_{(\theta'', r'')} \phi + \frac{n-k-1}{k} \mathbb{E}_{(\theta', r')} (1 - \phi) \right].$$

Moreover, there exists some  $i \in [n]$ , such that the two rank vectors  $r', r''$  satisfy  $\theta'_{r'_j} = \theta''_{r''_j}$  for all  $j \in [n] \setminus \{i\}$ . For the  $i$ th entry, we have  $\theta'_{r'_i} = -\Delta$  and  $\theta''_{r''_i} = 0$ . The reduction of the top- $k$  ranking problem to the testing problem (25) is the most important step in our lower bound analysis. A rigorous argument will be given in Appendix B.

The testing problem (25) can be roughly understood as to test whether the  $i$ th player belongs to the top- $k$  set or not. The two hypotheses receive different weights 1 and  $\frac{n-k-1}{k}$  because of the definition of the loss function  $H_k(\hat{r}, r^*)$ . The optimal procedure to (25) is given by the likelihood ratio test

$$\phi = \mathbb{I} \left\{ \frac{d\mathbb{P}_{(\theta', r')}}{d\mathbb{P}_{(\theta'', r'')}} \geq \frac{k}{n-k-1} \right\},$$

according to Neyman–Pearson lemma. Since the vectors  $\{\theta'_{r'_i}\}_{i \in [n]}$  and  $\{\theta''_{r''_i}\}_{i \in [n]}$  only differ at the  $i$ th entry, the likelihood ratio statistic only depends on  $\{\bar{y}_{ij}\}_{j \in [n] \setminus \{i\}}$  and  $\{A_{ij}\}_{j \in [n] \setminus \{i\}}$ . Therefore, the testing error (25) is relatively easy to quantify. A sharp lower bound can be obtained by a large deviation analysis.

**THEOREM 6.1.** *Assume  $\frac{np}{\log n} \rightarrow \infty$ ,  $\kappa \leq c_1$  for some constant  $c_1 > 0$ , and (16) holds for some arbitrarily small constant  $\epsilon > 0$ . Then there exists some  $\delta = o(1)$ , such that*

$$\begin{aligned} & \inf_{\hat{r}} \sup_{\substack{r^* \in \mathfrak{S}_n \\ \theta^* \in \Theta'(k, \Delta, \kappa)}} \mathbb{E}_{(\theta^*, r^*)} H_k(\hat{r}, r^*) \\ & \geq C \exp \left( -\frac{1}{2} \left( \frac{\sqrt{(1+\delta)\text{SNR}}}{2} - \frac{1}{\sqrt{(1+\delta)\text{SNR}}} \log \frac{n-k}{k} \right)_+^2 \right), \end{aligned}$$

for some constant  $C > 0$  only depending on  $c_1$  and  $\epsilon$ .

**7. Local error rates.** So far, our study of the top- $k$  ranking problem has been conducted under the minimax decision-theoretic framework laid out in Section 2. The upper and lower bounds for the MLE and the spectral method are established uniformly over the parameter space  $\Theta(k, \Delta, \kappa)$ . To complement the minimax results, in this section, we present local error rates for the MLE and the spectral method, which leads to a refined comparison between the two popular methods.

**7.1. Local error rate for the MLE.** To analyze the statistical property of the MLE for each individual  $\theta$ , we first need to generalize the effective variance (11). For any  $\theta \in \mathbb{R}^n$  and any  $i \in [n]$ , define

$$(26) \quad V_i(\theta) = \frac{n}{\sum_{j=1}^n \psi'(\theta_i - \theta_j)}.$$

With the help of  $V_i(\theta)$ , for any subset  $S \subset [n]$ , we also define

$$(27) \quad R_1(S, \theta, t, \delta) = \sum_{i \in S} \exp \left( -\frac{(1-\delta)(\theta_i - t)_+^2 npL}{2V_i(\theta)} \right),$$

$$(28) \quad R_2(S, \theta, t, \delta) = \sum_{i \in S} \exp \left( -\frac{(1-\delta)(t - \theta_i)_+^2 npL}{2V_i(\theta)} \right).$$

THEOREM 7.1. Assume  $\frac{np}{\log n} \rightarrow \infty, \kappa \leq c_1$  for some constant  $c_1 > 0$ . Then, for the rank vector  $\hat{r}$  that is induced by the MLE (6), any small constant  $0 < \delta < 0.1$ , any  $r^* \in \mathfrak{S}_n$  and any  $\theta^* \in \Theta(k, 0, \kappa)$ , we have

$$(29) \quad \mathbb{E}_{(\theta^*, r^*)} H_k(\hat{r}, r^*) \leq C_1 \left( \inf_t \frac{R_1([k], \theta^*, t, \delta) + R_2([n] \setminus [k], \theta^*, t, \delta)}{k} + n^{-3} \right),$$

where  $C_1 > 0$  is a constant only depending on  $c_1$  and  $\delta$ . Moreover, we also have

$$(30) \quad \mathbb{E}_{(\theta^*, r^*)} H_k(\hat{r}, r^*) \geq C_2 \left( \inf_t \frac{R_1([k], \theta^*, t, -\delta) + R_2([n] \setminus [k], \theta^*, t, -\delta)}{k} \right),$$

for some constant  $C_2 > 0$  only depending on  $c_1$  and  $\delta$ , if we additionally assume that  $\inf_t (R_1([k], \theta^*, t, -\delta) + R_2([n] \setminus [k], \theta^*, t, -\delta)) \rightarrow \infty$ .

Theorem 7.1 gives matching upper and lower bounds for the error of the MLE for each individual  $\theta^* \in \Theta(k, 0, \kappa)$  and  $r^* \in \mathfrak{S}_n$ , except for the additional  $n^{-3}$  term and an arbitrarily small  $\delta$ . We remark that the  $n^{-3}$  term in the upper bound can be replaced by  $n^{-C}$  for an arbitrarily large constant  $C$ . The upper bound (29) can be viewed as an extension of Theorem 3.2, though the  $\delta$  in Theorem 3.2 is allowed to vanish because of the less general setting. The lower bound (30) requires an extra condition  $\inf_t (R_1([k], \theta^*, t, -\delta) + R_2([n] \setminus [k], \theta^*, t, -\delta)) \rightarrow \infty$ , which implies the error rate is of higher order than  $O(k^{-1})$ . It plays the same role as the condition (20) in Theorem 4.3. This assumption covers most interesting partial recovery cases, since  $O(k^{-1})$  is already the error rate of exact recovery.

Let  $t^*$  be a minimizer of the right-hand side of (29) or (30). Then we can interpret  $\sum_{i=1}^k \exp(-\frac{(\theta_i^* - t^*)^2_{+} npL}{2V_i(\theta^*)})$  as the order of the number of top  $k$  players that are ranked among the bottom group, and  $\sum_{i=k+1}^n \exp(-\frac{(t^* - \theta_i^*)^2_{+} npL}{2V_i(\theta^*)})$  as the order of the number of bottom  $n - k$  players that are ranked in the top group.

A careful reader may notice that the error rate in Theorem 7.1 does not have a clear dependence on the signal gap  $\theta_k^* - \theta_{k+1}^*$ . This is because the current error rate depends on  $\theta^*$  more explicitly rather than just the difference between  $\theta_k^*$  and  $\theta_{k+1}^*$ . Even when  $\theta_k^* = \theta_{k+1}^*$ , it is still possible that the right-hand side of (29) converges to zero as long as the majority of  $\{\theta_i^*\}_{1 \leq i \leq k}$  are separated from most of  $\{\theta_i^*\}_{k+1 \leq i \leq n}$ .

7.2. *Local error rate for the spectral method.* To present a similar local error rate for the spectral method, we also need to generalize the effective variance (17). For any  $\theta \in \mathbb{R}^n$  and any  $i \in [n]$ , define

$$(31) \quad \bar{V}_i(\theta) = \frac{n \sum_{j=1}^n \psi'(\theta_i - \theta_j)(1 + e^{\theta_j - \theta_i})^2}{(\sum_{j=1}^n \psi(\theta_j - \theta_i))^2}.$$

We also introduce two quantities similar to (27) and (28),

$$\begin{aligned} \bar{R}_1(S, \theta, t, \delta) &= \sum_{i \in S} \exp\left(-\frac{(1 - \delta)(\theta_i - t)^2_{+} npL}{2\bar{V}_i(\theta)}\right), \\ \bar{R}_2(S, \theta, t, \delta) &= \sum_{i \in S} \exp\left(-\frac{(1 - \delta)(t - \theta_i)^2_{+} npL}{2\bar{V}_i(\theta)}\right). \end{aligned}$$

THEOREM 7.2. Assume  $\frac{np}{\log n} \rightarrow \infty, \kappa \leq c_1$  for some constant  $c_1 > 0$ . Then, for the rank vector  $\hat{r}$  that is induced by the stationary distribution of the Markov chain (7), any small

constant  $0 < \delta < 0.1$ , any  $r^* \in \mathfrak{S}_n$  and any  $\theta^* \in \Theta(k, 0, \kappa)$ , we have

$$(32) \quad \mathbb{E}_{(\theta^*, r^*)} H_k(\hat{r}, r^*) \leq C_1 \left( \inf_t \frac{\bar{R}_1([k], \theta^*, t, \delta) + \bar{R}_2([n] \setminus [k], \theta^*, t, \delta)}{k} + n^{-3} \right),$$

where  $C_1 > 0$  is a constant only depending on  $c_1$  and  $\delta$ . Moreover, we also have

$$(33) \quad \mathbb{E}_{(\theta^*, r^*)} H_k(\hat{r}, r^*) \geq C_2 \left( \inf_t \frac{\bar{R}_1([k], \theta^*, t, -\delta) + \bar{R}_2([n] \setminus [k], \theta^*, t, -\delta)}{k} \right),$$

for some constant  $C_2 > 0$  only depending on  $c_1$  and  $\delta$ , if we additionally assume that  $\inf_t (\bar{R}_1([k], \theta^*, t, -\delta) + \bar{R}_2([n] \setminus [k], \theta^*, t, -\delta)) \rightarrow \infty$ .

Similar to Theorem 7.1, Theorem 7.2 also gives matching upper and lower bounds for the error of the spectral method for each individual  $\theta^* \in \Theta(k, 0, \kappa)$  and  $r^* \in \mathfrak{S}_n$ .

Let us remark that the results of Theorem 7.1 and Theorem 7.2 can be further extended beyond the setting of the Erdős–Rényi graph and exactly  $L$  comparisons on each edge. To be specific, we can consider a random graph  $A_{ij} \sim \text{Bernoulli}(p_{ij})$  independently for all  $1 \leq i < j \leq n$ . For each edge, we observe  $L_{ij}$  independent games. Then, as long as  $\max_{ij} p_{ij} \leq C \min_{ij} p_{ij}$  and  $\max_{ij} L_{ij} \leq C \min_{ij} L_{ij}$  hold for some constant  $C > 0$ , the results of Theorem 7.1 and Theorem 7.2 continue to hold with  $\frac{V_i(\theta)}{npL}$  and  $\frac{\bar{V}_i(\theta)}{npL}$  replaced by  $\frac{\sum_{j \in [n] \setminus \{i\}} \frac{p_{ij}}{L_{ij}} \psi(\theta_i - \theta_j) \psi(\theta_j - \theta_i)}{(\sum_{j \in [n] \setminus \{i\}} p_{ij} \psi(\theta_i - \theta_j) \psi(\theta_j - \theta_i))^2}$  and  $\frac{\sum_{j \in [n] \setminus \{i\}} \frac{p_{ij}}{L_{ij}} \psi(\theta_i - \theta_j) \psi(\theta_j - \theta_i) (1 + e^{\theta_j - \theta_i})^2}{(\sum_{j \in [n] \setminus \{i\}} p_{ij} \psi(\theta_j - \theta_i))^2}$ , respectively.

**7.3. Comparison of the two methods for each  $\theta^*$ .** Theorem 7.1 and Theorem 7.2 allow us to give a refined comparison between the MLE and the spectral method. By ignoring the  $n^{-3}$  term in the upper bounds and the  $\delta$  in each exponent, we can write the error rates of the MLE and the spectral method as

$$(34) \quad \inf_t \frac{1}{k} \left[ \sum_{i=1}^k \exp\left(-\frac{(\theta_i^* - t)_+^2 npL}{2V_i(\theta^*)}\right) + \sum_{i=k+1}^n \exp\left(-\frac{(t - \theta_i^*)^2 npL}{2V_i(\theta^*)}\right) \right],$$

and

$$(35) \quad \inf_t \frac{1}{k} \left[ \sum_{i=1}^k \exp\left(-\frac{(\theta_i^* - t)_+^2 npL}{2\bar{V}_i(\theta^*)}\right) + \sum_{i=k+1}^n \exp\left(-\frac{(t - \theta_i^*)^2 npL}{2\bar{V}_i(\theta^*)}\right) \right].$$

It is clear that the only difference between (34) and (35) lies in the difference of the variance functions (26) and (31), whose comparison is given by the following lemma.

**LEMMA 7.1.** *For any  $\theta^* \in \mathbb{R}$  and any  $i \in [n]$ , we have  $V_i(\theta^*) \leq \bar{V}_i(\theta^*)$ . The equality holds if and only if  $\theta_1^* = \dots = \theta_n^*$ .*

**PROOF.** Notice the following chain of equalities and inequality:

$$\begin{aligned} V_i(\theta^*) &= \frac{n}{\sum_{j \in [n]} \psi'(\theta_j^* - \theta_i^*)} \\ &= \frac{n(\sum_{j \in [n]} e^{\theta_j^* - \theta_i^*})}{(\sum_{j \in [n]} \psi'(\theta_j^* - \theta_i^*))(\sum_{j \in [n]} e^{\theta_j^* - \theta_i^*})} \\ (36) \quad &\leq \frac{n(\sum_{j \in [n]} \psi'(\theta_j^* - \theta_i^*)(1 + e^{\theta_j^* - \theta_i^*})^2)}{(\sum_{j \in [n]} \psi(\theta_j^* - \theta_i^*))^2} \\ &= \bar{V}_i(\theta^*), \end{aligned}$$

where (36) is by the Cauchy–Schwarz inequality on the denominator. According to the equality condition of the Cauchy–Schwarz inequality, we know that  $V_i(\theta^*) = \overline{V}_i(\theta^*)$  only when  $\theta_1^* = \cdots = \theta_n^*$ .  $\square$

To close this section, we discuss two special cases of  $\theta^*$ , under which the error rates recover the results of Theorem 3.2, Theorem 4.1 and Corollary 5.1.

**EXAMPLE 7.1.** According to the proof of Theorem 4.3 and the construction discussed in Section 6, the least favorable  $\theta^* \in \Theta(k, \Delta, \kappa)$  takes the following form:  $\theta_i^* = \kappa_1$  for all  $1 \leq i \leq k - \rho k$ ,  $\theta_i^* = 0$  for  $k - \rho k < i \leq k$ ,  $\theta_i^* = -\Delta$  for  $k < i \leq k + \rho(n - k)$  and  $\theta_i^* = -\kappa_2$  for  $k + \rho(n - k) < i \leq n$ . Here,  $\kappa_1$  and  $\kappa_2$  are maximizers of either (11) for the MLE or (17) for the spectral method, and  $\rho$  is a sufficiently small constant. For this  $\theta^*$ , the formulas (34) and (35) recover the minimax rates obtained in Theorem 3.2 and Theorem 4.1.

**EXAMPLE 7.2.** Another interesting  $\theta^*$  is the two-piece model  $\theta^* \in \Theta(k, \Delta, \Delta)$ . By the translational invariance of the variance functions, we can consider  $\theta_i^* = \Delta$  for all  $1 \leq i \leq k$  and  $\theta_i^* = 0$  for all  $k < i \leq n$ . We discuss the consequence of this choice of  $\theta^*$  under two situations. First, consider  $\Delta = o(1)$ , and one can check that  $V_i(\theta^*) = (1 + o(1))4$  and  $\overline{V}_i(\theta^*) = (1 + o(1))4$  for all  $i \in [n]$ , which implies the equivalence of error rates of the MLE and the spectral method. Second, consider  $\Delta$  lower bounded by some constant. In this case, both the formulas (34) and (35) are  $o(k^{-1})$ , which implies both the MLE and the spectral method achieve exact recovery with high probability. As shown in Corollary 5.1, the spectral method is actually optimal for  $\theta^* \in \Theta(k, \Delta, \Delta)$ . We are therefore able to give a theoretical justification of the numerical experiment of [7].

**8. Analysis of the MLE.** In this section, we analyze the MLE (6), and prove Theorem 3.1, Theorem 3.2 and Theorem 3.3. Since the BTL model (3) is invariant to a shift of the model parameter, we can assume  $\mathbb{1}_n^T \theta^* = 0$  without loss of generality. For simplicity of notation, we also assume  $r_i^* = i$  for each  $i \in [n]$ , and thus we have  $\theta_{r_i^*}^* = \theta_i^*$ . Recall the convention of notation that  $A_{ij} = A_{ji}$  and  $\bar{y}_{ij} = 1 - \bar{y}_{ji}$  for any  $i < j$ . We also set  $A_{ii} = 0$  for all  $i \in [n]$ . Throughout the analysis, we will repeatedly use the properties that both  $\psi(t)$  and  $\psi'(t)$  are bounded continuous functions with bounded Lipschitz constants.

The section is organized as follows. We will first give a brief overview of the techniques and the main steps of the analysis in Section 8.1. We then present a few technical lemmas in Section 8.2. In Section 8.3, we establish an important result on the  $\ell_\infty$  bound of the MLE. Theorem 3.1 will be proved in Section 8.4. Finally, we prove Theorem 3.2 and Theorem 3.3 in Section 8.5.

**8.1. Overview of the techniques.** A major difficulty of analyzing the MLE is to control the spectrum of the Hessian matrix of the negative log-likelihood function. Recall the definition of  $\ell_n(\theta)$  in (5). Its Hessian  $\nabla^2 \ell_n(\theta) = H(\theta) \in \mathbb{R}^{n \times n}$  is given by the formula

$$H_{ij}(\theta) = \begin{cases} \sum_{l \in [n] \setminus \{i\}} A_{il} \psi'(\theta_i - \theta_l), & i = j, \\ -A_{ij} \psi'(\theta_i - \theta_j), & i \neq j. \end{cases}$$

It can be viewed as the Laplacian of the weighted random graph  $\{\psi'(\theta_i - \theta_j)A_{ij}\}$ . For  $\theta$  that satisfies  $\max_{i < j} |\theta_i - \theta_j| = O(1)$ , the spectrum of  $H(\theta)$  can be well controlled via some standard random matrix tool [24]. The property  $\max_{i < j} |\theta_i - \theta_j| = O(1)$  certainly holds for  $\theta^* \in \Theta(k, \Delta, \kappa)$ . However, when analyzing the Taylor expansion of  $\ell_n(\theta)$ , we actually need to understand  $H(\theta)$  for  $\theta$  that is a convex combination between  $\hat{\theta}$  and  $\theta^*$ . Since the MLE

is defined without any constraint or regularization, there is no such control for  $\hat{\theta}$ . Our first step is to establish the following proposition that shows  $\|\hat{\theta} - \theta^*\|_\infty$  is bounded with high probability even though the MLE has no constraint or regularization.

PROPOSITION 8.1. *Under the setting of Theorem 3.1, we have*

$$(37) \quad \|\hat{\theta} - \theta^*\|_\infty \leq 5,$$

with probability at least  $1 - O(n^{-7})$ .

The proof of Proposition 8.1 borrows strength from the property of a regularized MLE. Recall the definition of  $\hat{\theta}_\lambda$  in (10). This is the version of MLE that has been analyzed by [7]. We will choose  $\lambda = n^{-1}$  in order that  $\hat{\theta}_\lambda$  is close to  $\hat{\theta}$ . Following the techniques in [7], we can first show  $\|\hat{\theta}_\lambda - \theta^*\|_\infty \leq 4$  with high probability. The presence of the penalty in (10) is crucial for the result  $\|\hat{\theta}_\lambda - \theta^*\|_\infty \leq 4$  to be established. Next, we have an argument to show that the two estimators  $\hat{\theta}_\lambda$  and  $\hat{\theta}$  are sufficiently close. This leads to the bound (37). A detailed proof of Proposition 8.1 will be given in Section 8.3.

The result of Proposition 8.1 is arguably the most important step in the analysis of the MLE. It directly leads to the control of the spectrum of  $H(\theta)$ . Then the first bound (8) of Theorem 3.1 can be obtained by a Taylor expansion of the objective function  $\ell_n(\theta)$ . The second bound (9) of Theorem 3.1 and Theorem 3.2 requires an entrywise analysis of  $\hat{\theta}$ , and is therefore more complicated. We need to take advantage of the powerful leave-one-out argument in [7]. The intuition of the leave-one-out technique has been thoroughly discussed in [7], and we do not repeat it here. We would like to emphasize that our version of the leave-one-out argument is in fact different from the form introduced in [7]. We do not need to combine the leave-one-out argument with a gradient descent analysis as in [7]. This helps us to avoid the extra technical condition  $\log L = O(\log n)$  in [7] when proving the theorems.

8.2. *Some technical lemmas.* Let us present a few technical lemmas that facilitate our analysis of the MLE. The first two lemmas are concentration properties of the random graph  $A \sim \mathcal{G}(n, p)$ . We define  $\mathcal{L}_A = D - A$  to be the graph Laplacian of  $A$ , where  $D$  is a diagonal matrix whose entries are given by  $D_{ii} = \sum_{j \in [n] \setminus \{i\}} A_{ij}$ .

LEMMA 8.1. *Assume  $p \geq \frac{c_0 \log n}{n}$  for some sufficiently large  $c_0 > 0$ . We then have*

$$\frac{1}{2}np \leq \min_{i \in [n]} \sum_{j \in [n] \setminus \{i\}} A_{ij} \leq \max_{i \in [n]} \sum_{j \in [n] \setminus \{i\}} A_{ij} \leq 2np,$$

and

$$\lambda_{\min, \perp}(\mathcal{L}_A) = \min_{u \neq 0: \mathbf{1}_n^T u = 0} \frac{u^T \mathcal{L}_A u}{\|u\|^2} \geq \frac{np}{2},$$

$$\lambda_{\max}(\mathcal{L}_A) = \max_{u \neq 0} \frac{u^T \mathcal{L}_A u}{\|u\|^2} \leq 2np$$

with probability at least  $1 - O(n^{-10})$ .

LEMMA 8.2. *Assume  $p \geq \frac{c_0 \log n}{n}$  for some sufficiently large  $c_0 > 0$ . For any fixed  $\{w_{ij}\}$ , we have*

$$\max_{i \in [n]} \sum_{j \in [n] \setminus \{i\}} w_{ij}^2 (A_{ij} - p)^2 \leq Cnp \max_{i, j \in [n]} |w_{ij}|^2,$$

and

$$\max_{i \in [n]} \left( \sum_{j \in [n] \setminus \{i\}} w_{ij} (A_{ij} - p) \right)^2 \leq C(\log n)^2 \max_{i, j \in [n]} |w_{ij}|^2 + Cp \log n \max_{i \in [n]} \sum_{j \in [n]} w_{ij}^2,$$

for some constant  $C > 0$  with probability at least  $1 - O(n^{-10})$ .

With  $\lambda_{\min, \perp}(\mathcal{L}_A)$  shown to be well behaved, the next lemma establishes a similar control for  $\lambda_{\min, \perp}(H(\theta))$ .

LEMMA 8.3. Assume  $p \geq \frac{c_0 \log n}{n}$  for some sufficiently large  $c_0 > 0$ . For any  $\theta \in \mathbb{R}^n$  that satisfies  $\max_{i \in [n]} \theta_i - \min_{i \in [n]} \theta_i \leq M$ , we have

$$\lambda_{\min, \perp}(H(\theta)) \geq \frac{1}{8} e^{-Mnp},$$

with probability at least  $1 - O(n^{-10})$ .

Finally, we need a few concentration inequalities.

LEMMA 8.4. Assume  $\kappa = O(1)$  and  $p \geq \frac{c_0 \log n}{n}$  for some sufficiently large  $c_0 > 0$ . Then we have

$$\begin{aligned} \sum_{i=1}^n \left( \sum_{j \in [n] \setminus \{i\}} A_{ij} (\bar{y}_{ij} - \psi(\theta_i^* - \theta_j^*)) \right)^2 &\leq C \frac{n^2 p}{L}, \\ \max_{i \in [n]} \left( \sum_{j \in [n] \setminus \{i\}} A_{ij} (\bar{y}_{ij} - \psi(\theta_i^* - \theta_j^*)) \right)^2 &\leq C \frac{np \log n}{L}, \\ \max_{i \in [n]} \sum_{j \in [n] \setminus \{i\}} A_{ij} (\bar{y}_{ij} - \psi(\theta_i^* - \theta_j^*))^2 &\leq C \frac{np}{L}, \end{aligned}$$

for some constant  $C > 0$  with probability at least  $1 - O(n^{-10})$  uniformly over all  $\theta^* \in \Theta(k, 0, \kappa)$ .

The proofs of the four lemmas above will be given in Appendix E.

8.3. *Proof of Proposition 8.1.* As we have outlined in Section 8.1, the main argument to bound  $\|\hat{\theta} - \theta^*\|_\infty$  is to first derive a bound for  $\|\hat{\theta}_\lambda - \theta^*\|_\infty$ , where  $\hat{\theta}_\lambda$  is the penalized MLE defined in (10). Then we only need to show  $\hat{\theta}_\lambda$  and  $\hat{\theta}$  are close with  $\lambda$  as small as  $\lambda = n^{-1}$ . We first state a lemma that bounds  $\|\hat{\theta}_\lambda - \theta^*\|_\infty$ .

LEMMA 8.5. Under the setting of Theorem 3.1, for the estimator  $\hat{\theta}_\lambda$  with  $\lambda = n^{-1}$ , we have

$$\|\hat{\theta}_\lambda - \theta^*\|_\infty \leq 4,$$

with probability at least  $1 - O(n^{-7})$ .

We first prove Proposition 8.1 with the help of Lemma 8.5. The proof of Lemma 8.5 largely follows the arguments in [7] that analyze the regularized MLE. Since we only need to show  $\|\hat{\theta}_\lambda - \theta^*\|_\infty \leq 4$  rather than the optimal rate, the condition on  $L$  imposed by [7] is not needed anymore. This requires a few minor changes in the proof of [7]. We include the detailed proof of Lemma 8.5 in Appendix D in the Supplementary Material to be self-contained.

PROOF OF PROPOSITION 8.1. Define a constraint MLE as

$$(38) \quad \hat{\theta}^{\text{con}} = \underset{\mathbb{1}_n^T \theta = 0: \|\theta - \theta^*\|_\infty \leq 5}{\operatorname{argmin}} \ell_n(\theta).$$

By Lemma 8.5,  $\hat{\theta}_\lambda$  is feasible for the constraint of (38). We then have

$$(39) \quad \ell_n(\hat{\theta}_\lambda) \geq \ell_n(\hat{\theta}^{\text{con}}).$$

We apply Taylor expansion, and obtain

$$\ell_n(\hat{\theta}^{\text{con}}) = \ell_n(\hat{\theta}_\lambda) + (\hat{\theta}^{\text{con}} - \hat{\theta}_\lambda)^T \nabla \ell_n(\hat{\theta}_\lambda) + \frac{1}{2}(\hat{\theta}_\lambda - \hat{\theta}^{\text{con}})^T H(\xi)(\hat{\theta}_\lambda - \hat{\theta}^{\text{con}}),$$

where  $\xi$  is a convex combination of  $\hat{\theta}^{\text{con}}$  and  $\hat{\theta}_\lambda$ . By Lemma 8.5, we know that  $\|\hat{\theta}_\lambda - \theta^*\|_\infty \leq 4$ . We also have  $\|\hat{\theta}^{\text{con}} - \theta^*\|_\infty \leq 5$  by the definition of  $\hat{\theta}^{\text{con}}$ . Thus,  $\|\xi - \theta^*\|_\infty \leq 5$ . By Lemma 8.3, we get the lower bound

$$(40) \quad \ell_n(\hat{\theta}^{\text{con}}) \geq \ell_n(\hat{\theta}_\lambda) + (\hat{\theta}^{\text{con}} - \hat{\theta}_\lambda)^T \nabla \ell_n(\hat{\theta}_\lambda) + c_1 n p \|\hat{\theta}^{\text{con}} - \hat{\theta}_\lambda\|^2,$$

for some constant  $c_1 > 0$ . By (39) and (40), we have

$$\|\hat{\theta}^{\text{con}} - \hat{\theta}_\lambda\|^2 \leq \frac{|(\hat{\theta}^{\text{con}} - \hat{\theta}_\lambda)^T \nabla \ell_n(\hat{\theta}_\lambda)|}{c_1 n p}.$$

By the Cauchy–Schwarz inequality and the fact that  $\nabla \ell_n(\hat{\theta}_\lambda) + \lambda \hat{\theta}_\lambda = 0$ , we have

$$\|\hat{\theta}^{\text{con}} - \hat{\theta}_\lambda\|^2 \leq \frac{\|\nabla \ell_n(\hat{\theta}_\lambda)\|^2}{(c_1 n p)^2} = \frac{\lambda^2 \|\hat{\theta}_\lambda\|^2}{(c_1 n p)^2} \lesssim \frac{n \lambda^2}{(c_1 n p)^2} \lesssim n^{-1}.$$

Finally, since

$$\|\hat{\theta}^{\text{con}} - \theta^*\|_\infty \leq \|\hat{\theta}_\lambda - \theta^*\|_\infty + \|\hat{\theta}^{\text{con}} - \hat{\theta}_\lambda\| \leq 4 + \frac{c_2}{\sqrt{n}} \leq \frac{9}{2},$$

the minimizer of (38) is in the interior of the constraint. By the convexity of (38), we have  $\hat{\theta}^{\text{con}} = \hat{\theta}$ , and thus the desired conclusion  $\|\hat{\theta} - \theta^*\|_\infty \leq 5$  is obtained.  $\square$

8.4. *Proof of Theorem 3.1.* We give separate proofs for the conclusions (8) and (9) in this section.

PROOF OF (8) OF THEOREM 3.1. By the definition of  $\hat{\theta}$ , we have  $\ell_n(\theta^*) \geq \ell_n(\hat{\theta})$ . We then apply Taylor expansion and obtain

$$\ell_n(\hat{\theta}) = \ell_n(\theta^*) + (\hat{\theta} - \theta^*)^T \nabla \ell_n(\theta^*) + \frac{1}{2}(\hat{\theta} - \theta^*)^T H(\xi)(\hat{\theta} - \theta^*),$$

where  $\xi$  is a convex combination of  $\hat{\theta}$  and  $\theta^*$ . By Proposition 8.1, we have  $\|\hat{\theta} - \theta^*\|_\infty \leq 5$ , which implies  $\|\xi - \theta^*\|_\infty \leq 5$ . Thus, we can apply Lemma 8.3 and get  $\frac{1}{2}(\hat{\theta} - \theta^*)^T H(\xi)(\hat{\theta} - \theta^*) \geq c_1 n p \|\hat{\theta} - \theta^*\|^2$  for some constant  $c_1 > 0$ . Together with  $\ell_n(\theta^*) \geq \ell_n(\hat{\theta})$  and a Cauchy–Schwarz inequality, we have  $\|\hat{\theta} - \theta^*\|^2 \leq \frac{\|\nabla \ell_n(\theta^*)\|^2}{(c_1 n p)^2}$ . Use (S128) and Lemma 8.4, we obtain the desired conclusion that  $\|\hat{\theta} - \theta^*\|^2 \lesssim \frac{1}{L p}$ .  $\square$

The proof of (9) is more involved. It is based on a leave-one-out argument that is very different from the one used in [7]. Let us decompose the objective function  $\ell_n(\theta)$  as

$$(41) \quad \ell_n(\theta) = \ell_n^{(-m)}(\theta_{-m}) + \ell_n^{(m)}(\theta_m | \theta_{-m}),$$

where we use  $\theta_m \in \mathbb{R}$  for the  $m$ th entry of  $\theta$  and  $\theta_{-m} \in \mathbb{R}^{n-1}$  for the remaining entries. The two functions in (41) are defined as

$$\begin{aligned}\ell_n^{(-m)}(\theta_{-m}) &= \sum_{1 \leq i < j \leq n: i, j \neq m} A_{ij} \left[ \bar{y}_{ij} \log \frac{1}{\psi(\theta_i - \theta_j)} + (1 - \bar{y}_{ij}) \log \frac{1}{1 - \psi(\theta_i - \theta_j)} \right], \\ \ell_n^{(m)}(\theta_m | \theta_{-m}) &= \sum_{j \in [n] \setminus \{m\}} A_{mj} \left[ \bar{y}_{mj} \log \frac{1}{\psi(\theta_m - \theta_j)} + (1 - \bar{y}_{mj}) \log \frac{1}{1 - \psi(\theta_m - \theta_j)} \right].\end{aligned}$$

Define

$$(42) \quad \theta_{-m}^{(m)} = \underset{\theta_{-m}: \|\theta_{-m} - \theta_{-m}^*\|_\infty \leq 5}{\operatorname{argmin}} \ell_n^{(-m)}(\theta_{-m}).$$

We first present an  $\ell_2$  norm bound for  $\theta_{-m}^{(m)}$ . We also use  $H^{(-m)}(\theta_{-m})$  for the Hessian matrix  $\nabla^2 \ell_n^{(-m)}(\theta_{-m})$ .

LEMMA 8.6. *Under the setting of Theorem 3.1, there exists some constant  $C > 0$  such that*

$$\max_{m \in [n]} \|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbf{1}_{n-1}\|^2 \leq C \frac{1}{pL},$$

with probability at least  $1 - O(n^{-9})$ , where  $a_m = \operatorname{ave}(\theta_{-m}^{(m)} - \theta_{-m}^*)$ .

PROOF. The proof is very similar to that of (8), since  $\theta_{-m}^{(m)}$  can be thought of as a constrained MLE on a subset of the data. By the definition of  $\theta_{-m}^{(m)}$ , we have

$$\begin{aligned}\ell_n^{(-m)}(\theta_{-m}^*) &\geq \ell_n^{(-m)}(\theta_{-m}^{(m)}) \\ &= \ell_n^{(-m)}(\theta_{-m}^*) + (\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbf{1}_{n-1})^T \nabla \ell_n^{(-m)}(\theta_{-m}^*) \\ &\quad + \frac{1}{2} (\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbf{1}_{n-1})^T H^{(-m)}(\xi) (\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbf{1}_{n-1}),\end{aligned}$$

where  $\xi$  is a convex combination of  $\theta_{-m}^{(m)}$  and  $\theta_{-m}^*$ . In the above Taylor expansion, we have also used the property that  $\ell_n^{(-m)}(\theta_{-m}) = \ell_n^{(-m)}(\theta_{-m} + c \mathbf{1}_{n-1})$ ,  $\nabla \ell_n^{(-m)}(\theta_{-m}) = \nabla \ell_n^{(-m)}(\theta_{-m} + c \mathbf{1}_{n-1})$  and  $H^{(-m)}(\theta_{-m}) = H^{(-m)}(\theta_{-m} + c \mathbf{1}_{n-1})$  for any  $c \in \mathbb{R}$ . Since  $\|\xi - \theta_{-m}^*\|_\infty \leq \|\theta_{-m}^{(m)} - \theta_{-m}^*\|_\infty \leq 5$ , we can apply Lemma 8.3 to the subset of the data, and obtain

$$\begin{aligned}&\frac{1}{2} (\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbf{1}_{n-1})^T H^{(-m)}(\xi) (\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbf{1}_{n-1}) \\ &\geq c_1 np \|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbf{1}_{n-1}\|^2,\end{aligned}$$

with probability at least  $1 - O(n^{-10})$  for some constant  $c_1 > 0$ . By the Cauchy–Schwarz inequality, we have

$$\|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbf{1}_{n-1}\|^2 \leq \frac{\|\nabla \ell_n^{(-m)}(\theta_{-m}^*)\|^2}{(c_1 np)^2}.$$

Apply (S128) and Lemma 8.4 to the subset of the data, and we obtain that  $\|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbf{1}_{n-1}\|^2 \leq C \frac{1}{pL}$  with probability at least  $1 - O(n^{-10})$ . Finally, a union bound argument leads to the desired result.  $\square$

With the help of Lemma 8.6, we are ready to prove (9).

PROOF OF (9) OF THEOREM 3.1. By Proposition 8.1, we have  $\|\widehat{\theta}_{-m} - \theta_{-m}^*\|_\infty \leq \|\widehat{\theta} - \theta^*\|_\infty \leq 5$ , and thus  $\widehat{\theta}_{-m}$  is feasible for the constraint of (42). By the definition of  $\theta_{-m}^{(m)}$ , we have

$$\begin{aligned} \ell_n^{(-m)}(\widehat{\theta}_{-m}) &\geq \ell_n^{(-m)}(\theta_{-m}^{(m)}) \\ &= \ell_n^{(-m)}(\widehat{\theta}_{-m}) + (\theta_{-m}^{(m)} - \widehat{\theta}_{-m} - \bar{a}_m \mathbf{1}_{n-1})^T \nabla \ell_n^{(-m)}(\widehat{\theta}_{-m}) \\ &\quad + \frac{1}{2}(\theta_{-m}^{(m)} - \widehat{\theta}_{-m} - \bar{a}_m \mathbf{1}_{n-1})^T H^{(-m)}(\xi)(\theta_{-m}^{(m)} - \widehat{\theta}_{-m} - \bar{a}_m \mathbf{1}_{n-1}), \end{aligned}$$

where  $\bar{a}_m = \text{ave}(\theta_{-m}^{(m)} - \widehat{\theta}_{-m})$  and  $\xi$  is a convex combination of  $\theta_{-m}^{(m)}$  and  $\widehat{\theta}_{-m}$ . Since both  $\theta_{-m}^{(m)}$  and  $\widehat{\theta}_{-m}$  satisfy the constraint of (42), we must have  $\|\xi - \theta_{-m}^*\|_\infty \leq 5$ . Then we can apply Lemma 8.3 to the subset of the data, and obtain

$$\begin{aligned} &\frac{1}{2}(\theta_{-m}^{(m)} - \widehat{\theta}_{-m} - \bar{a}_m \mathbf{1}_{n-1})^T H^{(-m)}(\xi)(\theta_{-m}^{(m)} - \widehat{\theta}_{-m} - \bar{a}_m \mathbf{1}_{n-1}) \\ &\geq c_1 np \|\theta_{-m}^{(m)} - \widehat{\theta}_{-m} - \bar{a}_m \mathbf{1}_{n-1}\|^2, \end{aligned}$$

for some constant  $c_1 > 0$ . By the Cauchy–Schwarz inequality, we have

$$\|\theta_{-m}^{(m)} - \widehat{\theta}_{-m} - \bar{a}_m \mathbf{1}_{n-1}\|^2 \leq \frac{\|\nabla \ell_n^{(-m)}(\widehat{\theta}_{-m})\|^2}{(c_1 np)^2}.$$

For each  $i \in [n] \setminus \{m\}$ , by the decomposition (41), we have

$$\frac{\partial}{\partial \theta_i} \ell_n^{(-m)}(\theta_{-m}) = \frac{\partial}{\partial \theta_i} \ell_n(\theta) - \frac{\partial}{\partial \theta_i} \ell_n^{(m)}(\theta_m | \theta_{-m}).$$

Since  $\nabla \ell_n(\widehat{\theta}) = 0$ , we have

$$\frac{\partial}{\partial \theta_i} \ell_n^{(-m)}(\theta_{-m})|_{\theta=\widehat{\theta}} = -\frac{\partial}{\partial \theta_i} \ell_n^{(m)}(\theta_m | \theta_{-m})|_{\theta=\widehat{\theta}} = -A_{mi}(\bar{y}_{mi} - \psi(\widehat{\theta}_m - \widehat{\theta}_i)).$$

We therefore have the bound,

$$\begin{aligned} \|\nabla \ell_n^{(-m)}(\widehat{\theta}_{-m})\|^2 &= \sum_{i \in [n] \setminus \{m\}} A_{mi}(\bar{y}_{mi} - \psi(\widehat{\theta}_m - \widehat{\theta}_i))^2 \\ &\leq 2 \sum_{i \in [n] \setminus \{m\}} A_{mi}(\bar{y}_{mi} - \psi(\theta_m^* - \theta_i^*))^2 \\ &\quad + 2 \sum_{i \in [n] \setminus \{m\}} A_{mi}(\psi(\theta_m^* - \theta_i^*) - \psi(\widehat{\theta}_m - \widehat{\theta}_i))^2 \\ &\leq 2 \sum_{i \in [n] \setminus \{m\}} A_{mi}(\bar{y}_{mi} - \psi(\theta_m^* - \theta_i^*))^2 + 2\|\widehat{\theta} - \theta^*\|_\infty^2 \sum_{i \in [n] \setminus \{m\}} A_{mi} \\ &\leq 2 \sum_{i \in [n] \setminus \{m\}} A_{mi}(\bar{y}_{mi} - \psi(\theta_m^* - \theta_i^*))^2 + 4np\|\widehat{\theta} - \theta^*\|_\infty^2, \end{aligned}$$

where the last inequality is by Lemma 8.1. This implies

$$\begin{aligned} \max_{m \in [n]} \|\theta_{-m}^{(m)} - \widehat{\theta}_{-m} - \bar{a}_m \mathbf{1}_{n-1}\|^2 &\leq \frac{\max_{m \in [n]} \sum_{i \in [n] \setminus \{m\}} A_{mi}(\bar{y}_{mi} - \psi(\theta_m^* - \theta_i^*))^2}{(c_1 np)^2/2} \\ &\quad + \frac{\|\widehat{\theta} - \theta^*\|_\infty^2}{c_1^2 np/4}. \end{aligned}$$

Since we need a bound for  $\max_{m \in [n]} \|\theta_{-m}^{(m)} - \widehat{\theta}_{-m} - a_m \mathbb{1}_{n-1}\|^2$ , we need to quantify the difference between  $a_m$  and  $\bar{a}_m$ . Recall that  $a_m = \text{ave}(\theta_{-m}^{(m)} - \theta_m^*)$ . Since  $\mathbb{1}_n^T \widehat{\theta} = \mathbb{1}_n^T \theta^* = 0$ , we have

$$\|a_m \mathbb{1}_{n-1} - \bar{a}_m \mathbb{1}_{n-1}\|^2 = (n-1)(\text{ave}(\widehat{\theta}_{-m} - \theta_m^*))^2 = \frac{(\widehat{\theta}_m - \theta_m^*)^2}{n-1} \leq \frac{\|\widehat{\theta} - \theta^*\|_\infty^2}{n-1}.$$

We then have

$$\begin{aligned} & \max_{m \in [n]} \|\theta_{-m}^{(m)} - \widehat{\theta}_{-m} - a_m \mathbb{1}_{n-1}\|^2 \\ (43) \quad & \leq C_1 \frac{\max_{m \in [n]} \sum_{i \in [n] \setminus \{m\}} A_{mi} (\bar{y}_{mi} - \psi(\theta_m^* - \theta_i^*))^2}{n^2 p^2} \\ & + C_1 \frac{\|\widehat{\theta} - \theta^*\|_\infty^2}{np}, \end{aligned}$$

for some constant  $C_1 > 0$ .

Next, let us derive a bound for  $\|\widehat{\theta} - \theta^*\|_\infty^2$  in terms of  $\max_{m \in [n]} \|\theta_{-m}^{(m)} - \widehat{\theta}_{-m} - a_m \mathbb{1}_{n-1}\|^2$ . We introduce the notation:

$$\begin{aligned} f^{(m)}(\theta_m | \theta_{-m}) &= \frac{\partial}{\partial \theta_m} \ell_n^{(m)}(\theta_m | \theta_{-m}) = - \sum_{i \in [n] \setminus \{m\}} A_{mi} (\bar{y}_{mi} - \psi(\theta_m - \theta_i)), \\ g^{(m)}(\theta_m | \theta_{-m}) &= \frac{\partial^2}{\partial \theta_m^2} \ell_n^{(m)}(\theta_m | \theta_{-m}) = \sum_{i \in [n] \setminus \{m\}} A_{mi} \psi(\theta_m - \theta_i) \psi(\theta_i - \theta_m). \end{aligned}$$

By the definition of  $\widehat{\theta}$ , we know that  $\ell_n(\widehat{\theta}) = \min_{\theta: \mathbb{1}_n^T \theta = 0} \ell_n(\theta)$ . Since  $\ell_n(\theta) = \ell_n(\theta + c \mathbb{1}_n)$  for any  $c \in \mathbb{R}$ , we also have  $\ell_n(\widehat{\theta}) = \min_{\theta} \ell_n(\theta)$ . This allows us to compare the value of the objective  $\ell_n(\theta)$  at  $\widehat{\theta}$  with any vector that is not necessarily centered. We then have

$$\ell_n^{(m)}(\theta_m^* | \widehat{\theta}_{-m}) + \ell_n^{(-m)}(\widehat{\theta}_{-m}) \geq \ell_n(\widehat{\theta}),$$

which implies

$$\begin{aligned} \ell_n^{(m)}(\theta_m^* | \widehat{\theta}_{-m}) &\geq \ell_n^{(m)}(\widehat{\theta}_m | \widehat{\theta}_{-m}) \\ &= \ell_n^{(m)}(\theta_m^* | \widehat{\theta}_{-m}) + (\widehat{\theta}_m - \theta_m^*) f^{(m)}(\theta_m^* | \widehat{\theta}_{-m}) + \frac{1}{2} (\widehat{\theta}_m - \theta_m^*)^2 g^{(m)}(\xi | \widehat{\theta}_{-m}), \end{aligned}$$

where  $\xi$  is a scalar between  $\theta_m^*$  and  $\widehat{\theta}_m$ . By Proposition 8.1,  $|\xi - \theta_m^*| \leq |\widehat{\theta}_m - \theta_m^*| \leq \|\widehat{\theta} - \theta^*\|_\infty \leq 5$ . Therefore, for any  $i \in [n] \setminus \{m\}$ ,  $|\xi - \widehat{\theta}_i| \leq |\xi - \theta_m^*| + |\theta_m^* - \theta_i^*| + |\widehat{\theta}_i - \theta_i^*| \leq 10 + \kappa$ . This implies  $\frac{1}{2} g^{(m)}(\xi | \widehat{\theta}_{-m}) \geq c_2 np$  for some constant  $c_2 > 0$  with the help of Lemma 8.1. We then have the bound

$$(44) \quad (\widehat{\theta}_m - \theta_m^*)^2 \leq \frac{|f^{(m)}(\theta_m^* | \widehat{\theta}_{-m})|^2}{(c_2 np)^2}.$$

We bound  $|f^{(m)}(\theta_m^* | \widehat{\theta}_{-m})|$  by

$$\begin{aligned} |f^{(m)}(\theta_m^* | \widehat{\theta}_{-m})| &= \left| \sum_{i \in [n] \setminus \{m\}} A_{mi} (\bar{y}_{mi} - \psi(\theta_m^* - \widehat{\theta}_i)) \right| \\ (45) \quad &\leq \left| \sum_{i \in [n] \setminus \{m\}} A_{mi} (\bar{y}_{mi} - \psi(\theta_m^* - \theta_i^*)) \right| \end{aligned}$$

$$(46) \quad + \left| \sum_{i \in [n] \setminus \{m\}} A_{mi} (\psi(\theta_m^* - \theta_i^*) - \psi(\theta_m^* - \theta_i^{(m)} + a_m)) \right|$$

$$(47) \quad + \left| \sum_{i \in [n] \setminus \{m\}} A_{mi} (\psi(\theta_m^* - \theta_i^{(m)} + a_m) - \psi(\theta_m^* - \hat{\theta}_i)) \right|.$$

We use Lemma 8.2 to bound (46). We have

$$(48) \quad \left| \sum_{i \in [n] \setminus \{m\}} A_{mi} (\psi(\theta_m^* - \theta_i^*) - \psi(\theta_m^* - \theta_i^{(m)} + a_m)) \right|$$

$$(49) \quad \leq p \left| \sum_{i \in [n] \setminus \{m\}} (\psi(\theta_m^* - \theta_i^*) - \psi(\theta_m^* - \theta_i^{(m)} + a_m)) \right|$$

$$(49) \quad + \left| \sum_{i \in [n] \setminus \{m\}} (A_{mi} - p)(\psi(\theta_m^* - \theta_i^*) - \psi(\theta_m^* - \theta_i^{(m)} + a_m)) \right|$$

$$\leq p\sqrt{n} \|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbb{1}_{n-1}\| + C_2 \log n \|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbb{1}_{n-1}\|_\infty$$

$$+ C_2 \sqrt{p \log n} \|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbb{1}_{n-1}\|$$

$$\leq (p\sqrt{n} + C_2 \sqrt{p \log n}) \|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbb{1}_{n-1}\| + C_2 \log n \|\hat{\theta} - \theta^*\|_\infty$$

$$+ C_2 \log n \|\theta_{-m}^{(m)} - \hat{\theta}_{-m} - a_m \mathbb{1}_{n-1}\|.$$

With the help of 8.1, we can also bound (47), and we get

$$(50) \quad \left| \sum_{i \in [n] \setminus \{m\}} A_{mi} (\psi(\theta_m^* - \theta_i^{(m)} + a_m) - \psi(\theta_m^* - \hat{\theta}_i)) \right|$$

$$\leq \sqrt{\sum_{i \in [n] \setminus \{m\}} A_{mi}} \|\theta_{-m}^{(m)} - \hat{\theta}_{-m} - a_m \mathbb{1}_{n-1}\|$$

$$\leq C_3 \sqrt{np} \|\theta_{-m}^{(m)} - \hat{\theta}_{-m} - a_m \mathbb{1}_{n-1}\|.$$

Plug the bounds into (44), and we have

$$\|\hat{\theta} - \theta^*\|_\infty \leq \frac{\max_{m \in [n]} |\sum_{i \in [n] \setminus \{m\}} A_{mi} (\bar{y}_{mi} - \psi(\theta_m^* - \theta_i^*))|}{c_2 np}$$

$$+ \frac{(p\sqrt{n} + C_2 \sqrt{p \log n}) \max_{m \in [n]} \|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbb{1}_{n-1}\|}{c_2 np}$$

$$+ \frac{(C_2 \log n + C_3 \sqrt{np}) \|\theta_{-m}^{(m)} - \hat{\theta}_{-m} - a_m \mathbb{1}_{n-1}\|}{c_2 np} + \frac{C_2 \log n \|\hat{\theta} - \theta^*\|_\infty}{c_2 np}.$$

Since  $np \geq c_0 \log n$  for some sufficiently large  $c_0$ , we obtain the bound

$$(51) \quad \|\hat{\theta} - \theta^*\|_\infty \leq C_4 \frac{\max_{m \in [n]} |\sum_{i \in [n] \setminus \{m\}} A_{mi} (\bar{y}_{mi} - \psi(\theta_m^* - \theta_i^*))|}{np}$$

$$+ C_4 \frac{p\sqrt{n} \max_{m \in [n]} \|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbb{1}_{n-1}\|}{np}$$

$$+ C_4 \frac{(\log n + \sqrt{np}) \|\theta_{-m}^{(m)} - \hat{\theta}_{-m} - a_m \mathbb{1}_{n-1}\|}{np}.$$

Let us plug the above bound into (43). Then, after some rearrangement, we obtain

$$\begin{aligned} \max_{m \in [n]} \|\theta_{-m}^{(m)} - \hat{\theta}_{-m} - a_m \mathbb{1}_{n-1}\| &\leq C_5 \frac{\max_{m \in [n]} \sqrt{\sum_{i \in [n] \setminus \{m\}} A_{mi} (\bar{y}_{mi} - \psi(\theta_m^* - \theta_i^*))^2}}{np} \\ &\quad + C_5 \frac{\max_{m \in [n]} |\sum_{i \in [n] \setminus \{m\}} A_{mi} (\bar{y}_{mi} - \psi(\theta_m^* - \theta_i^*))|}{np\sqrt{np}} \\ &\quad + C_5 \frac{\max_{m \in [n]} \|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbb{1}_{n-1}\|}{n\sqrt{p}}. \end{aligned}$$

By Lemma 8.4 and Lemma 8.6, we have

$$(52) \quad \max_{m \in [n]} \|\theta_{-m}^{(m)} - \hat{\theta}_{-m} - a_m \mathbb{1}_{n-1}\| \leq C_7 \sqrt{\frac{1}{npL}}.$$

Now we can plug the bound (52) back into (51), and together with Lemma 8.4 and Lemma 8.6, we have

$$(53) \quad \|\hat{\theta} - \theta^*\|_\infty \leq C_8 \sqrt{\frac{\log n}{npL}},$$

which is the desired conclusion. Tracking all the probabilistic events that we have used in the proof, we can conclude that both (52) and (53) hold with probability at least  $1 - O(n^{-7})$ .  $\square$

**8.5. Proofs of Theorem 3.2 and Theorem 3.3.** In the proof of (9), we have established the byproduct (52). This bound turns out to be extremely important for us to establish the result of Theorem 3.2. We therefore list it, together with its consequence, as a lemma.

**LEMMA 8.7.** *Under the setting of Theorem 3.1, there exists some constant  $C > 0$  such that*

$$\begin{aligned} \max_{m \in [n]} \|\theta_{-m}^{(m)} - \hat{\theta}_{-m} - a_m \mathbb{1}_{n-1}\|^2 &\leq C \frac{1}{npL}, \\ \max_{m \in [n]} \|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbb{1}_{n-1}\|_\infty^2 &\leq C \frac{\log n}{npL}, \end{aligned}$$

with probability at least  $1 - O(n^{-7})$ , where  $a_m = \text{ave}(\theta_{-m}^{(m)} - \theta_m^*)$  and  $\theta_{-m}^{(m)}$  is defined by (42).

**PROOF.** The first conclusion has been established in (52). The second conclusion is a consequence of the inequality

$$\max_{m \in [n]} \|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbb{1}_{n-1}\|_\infty^2 \leq 2 \max_{m \in [n]} \|\theta_{-m}^{(m)} - \hat{\theta}_{-m} - a_m \mathbb{1}_{n-1}\|^2 + 2 \|\hat{\theta} - \theta^*\|_\infty^2,$$

and (53).  $\square$

Now we are ready to prove Theorem 3.2.

**PROOF OF THEOREM 3.2.** When the error exponent is of constant order, the bound is also a constant, and the result already holds since  $H_k(\hat{r}, r^*) \leq 1$ . Therefore, we only need to consider the case when the error exponent tends to infinity. We first introduce some notation. Define

$$(54) \quad \eta = \frac{1}{2} - \frac{V(\kappa)}{(1 - \bar{\delta})\Delta^2 npL} \log \frac{n - k}{k},$$

where  $\bar{\delta} = o(1)$  is chosen such that  $\eta > 0$ . The specific choice of  $\bar{\delta}$  will be specified in the proof. Then let

$$(55) \quad \bar{\Delta}_i = \begin{cases} \min\left(\eta(\theta_k^* - \theta_{k+1}^*) + \theta_i^* - \theta_k^*, \left(\frac{\log n}{np}\right)^{1/4}\right), & 1 \leq i \leq k, \\ \min\left((1 - \eta)(\theta_k^* - \theta_{k+1}^*) + \theta_{k+1}^* - \theta_i^*, \left(\frac{\log n}{np}\right)^{1/4}\right), & k + 1 \leq i \leq n. \end{cases}$$

Since the diverging error exponent implies  $\text{SNR} \rightarrow \infty$ , we have  $\min_{i \in [n]} \bar{\Delta}_i^2 Lnp \rightarrow \infty$  and  $\max_{i \in [n]} \bar{\Delta}_i \rightarrow 0$ .

The proof involves several steps. In the first step, we need to derive a sharp probabilistic bound for  $|f^{(m)}(\theta_m^*|\hat{\theta}_{-m})|$ . In the proof of (9) of Theorem 3.1, we have shown that  $|f^{(m)}(\theta_m^*|\hat{\theta}_{-m})|$  can be bounded by the sum of (45), (47), (48) and (49). For (45), we can use Hoeffding's inequality and Lemma 8.1 and obtain the bound

$$(56) \quad \left| \sum_{i \in [n] \setminus \{m\}} A_{mi} (\bar{y}_{mi} - \psi(\theta_m^* - \theta_i^*)) \right| \leq C_1 \sqrt{\frac{x \sum_{i \in [n] \setminus \{m\}} A_{mi}}{L}} \leq C_2 \sqrt{\frac{xnp}{L}},$$

with probability at least  $1 - O(n^{-10}) - e^{-x}$ . Take  $x = \bar{\Delta}_m^{3/2} Lnp$ , and we have

$$(57) \quad \left| \sum_{i \in [n] \setminus \{m\}} A_{mi} (\bar{y}_{mi} - \psi(\theta_m^* - \theta_i^*)) \right| \leq C_2 \sqrt{\bar{\Delta}_m^{3/2} (np)^2},$$

with probability at least  $1 - O(n^{-10}) - e^{-\bar{\Delta}_m^{3/2} Lnp}$ . Since we have already shown (47) can be bounded by (50) with probability at least  $1 - O(n^{-10})$ , an application of Lemma 8.7 implies that

$$(58) \quad \left| \sum_{i \in [n] \setminus \{m\}} A_{mi} (\psi(\theta_m^* - \theta_i^{(m)} + a_m) - \psi(\theta_m^* - \hat{\theta}_i)) \right| \leq C_3 \sqrt{\frac{1}{L}},$$

with probability at least  $1 - O(n^{-7})$ . By the Cauchy-Schwarz inequality, we can bound (48) by  $p\sqrt{n}\|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbf{1}_{n-1}\|$ . With the help of Lemma 8.6, we have

$$(59) \quad p \left| \sum_{i \in [n] \setminus \{m\}} (\psi(\theta_m^* - \theta_i^*) - \psi(\theta_m^* - \theta_i^{(m)} + a_m)) \right| \leq C_4 \sqrt{\frac{np}{L}},$$

with probability at least  $1 - O(n^{-9})$ . For (49), we use Bernstein's inequality, and we have

$$(60) \quad \begin{aligned} & \left| \sum_{i \in [n] \setminus \{m\}} (A_{mi} - p)(\psi(\theta_m^* - \theta_i^*) - \psi(\theta_m^* - \theta_i^{(m)} + a_m)) \right| \\ & \leq C_5 \sqrt{px} \|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbf{1}_{n-1}\| + C_5 x \|\theta_{-m}^{(m)} - \theta_{-m}^* - a_m \mathbf{1}_{n-1}\|_\infty, \end{aligned}$$

with probability at least  $1 - e^{-x}$ . We choose  $x = \min(\bar{\Delta}_m^2 Lnp \frac{np}{\log n}, 7 \log n)$ . Then, with the help of Lemma 8.6 and Lemma 8.7, we have

$$(61) \quad \begin{aligned} & \left| \sum_{i \in [n] \setminus \{m\}} (A_{mi} - p)(\psi(\theta_m^* - \theta_i^*) - \psi(\theta_m^* - \theta_i^{(m)} + a_m)) \right| \\ & \leq C_6 \frac{1}{\sqrt{L}} \sqrt{\min\left(\bar{\Delta}_m^2 Lnp \frac{np}{\log n}, 7 \log n\right)} \\ & \quad + C_6 \sqrt{\frac{\log n}{npL}} \min\left(\bar{\Delta}_m^2 Lnp \frac{np}{\log n}, 7 \log n\right), \end{aligned}$$

with probability at least  $1 - O(n^{-7}) - \exp(-\bar{\Delta}_m^2 np L \frac{np}{\log n})$ . Combining the bounds (57)–(61), we obtain a bound for  $|f^{(m)}(\theta_m^*|\hat{\theta}_{-m})|$ . This also implies a bound for  $|\hat{\theta}_m - \theta_m^*|$  because of the inequality (44).

In the second step, we define

$$(62) \quad \bar{\theta}_m = \theta_m^* - \frac{f^{(m)}(\theta_m^*|\hat{\theta}_{-m})}{g^{(m)}(\theta_m^*|\hat{\theta}_{-m})}.$$

We need to show  $\bar{\theta}_m$  and  $\hat{\theta}_m$  are close. By Proposition 8.1,  $\|\hat{\theta} - \theta^*\|_\infty \leq 5$ , and thus  $g^{(m)}(\theta_m^*|\hat{\theta}_{-m}) \geq c_1 np$  for some constant  $c_1 > 0$ , so that we have the bound  $|\bar{\theta}_m - \theta_m^*| \leq \frac{|f^{(m)}(\theta_m^*|\hat{\theta}_{-m})|}{c_1 np}$ . In fact, given the inequality (44), we can choose  $c_1$  to be sufficiently small so that  $|\hat{\theta}_m - \theta_m^*| \leq \frac{|f^{(m)}(\theta_m^*|\hat{\theta}_{-m})|}{c_1 np}$  is also true. Therefore, we can express  $\bar{\theta}_m$  and  $\hat{\theta}_m$  as

$$\begin{aligned} \bar{\theta}_m &= \underset{|\theta_m - \theta_m^*| \leq \frac{|f^{(m)}(\theta_m^*|\hat{\theta}_{-m})|}{c_1 np}}{\operatorname{argmin}} \quad \bar{\ell}_n^{(m)}(\theta_m|\hat{\theta}_{-m}), \\ \hat{\theta}_m &= \underset{|\theta_m - \theta_m^*| \leq \frac{|f^{(m)}(\theta_m^*|\hat{\theta}_{-m})|}{c_1 np}}{\operatorname{argmin}} \quad \ell_n^{(m)}(\theta_m|\hat{\theta}_{-m}), \end{aligned}$$

where

$$\bar{\ell}_n^{(m)}(\theta_m|\hat{\theta}_{-m}) = \ell_n^{(m)}(\theta_m^*|\hat{\theta}_{-m}) + (\theta_m - \theta_m^*) f^{(m)}(\theta_m^*|\hat{\theta}_{-m}) + \frac{1}{2}(\theta_m - \theta_m^*)^2 g^{(m)}(\theta_m^*|\hat{\theta}_{-m}).$$

Recall the definition of  $\ell_n^{(m)}(\theta_m|\hat{\theta}_{-m})$  in (41) and the display afterwards. We will show  $\bar{\theta}_m$  and  $\hat{\theta}_m$  are close by bounding the difference between the two objective functions. By Taylor expansion, we have

$$|\ell_n^{(m)}(\theta_m|\hat{\theta}_{-m}) - \bar{\ell}_n^{(m)}(\theta_m|\hat{\theta}_{-m})| = \frac{1}{2}(\theta_m - \theta_m^*)^2 |g^{(m)}(\xi|\hat{\theta}_{-m}) - g^{(m)}(\theta_m^*|\hat{\theta}_{-m})|,$$

where  $\xi$  is a scalar between  $\theta_m$  and  $\theta_m^*$ . We then have

$$\begin{aligned} &|g^{(m)}(\xi|\hat{\theta}_{-m}) - g^{(m)}(\theta_m^*|\hat{\theta}_{-m})| \\ &= \left| \sum_{i \in [n] \setminus \{m\}} A_{mi} \psi(\xi - \hat{\theta}_i) \psi(\hat{\theta}_i - \xi) - \sum_{i \in [n] \setminus \{m\}} A_{mi} \psi(\theta_m^* - \hat{\theta}_i) \psi(\hat{\theta}_i - \theta_m^*) \right| \\ &\leq |\xi - \theta_m^*| \sum_{i \in [n] \setminus \{m\}} A_{mi} \\ &\leq C_7 |\theta_m - \theta_m^*| np, \end{aligned}$$

where the last inequality uses Lemma 8.1. Therefore, for any  $\theta_m$  that satisfies  $|\theta_m - \theta_m^*| \leq \frac{|f^{(m)}(\theta_m^*|\hat{\theta}_{-m})|}{c_1 np}$ , the difference between the two objective functions can be bounded by

$$|\ell_n^{(m)}(\theta_m|\hat{\theta}_{-m}) - \bar{\ell}_n^{(m)}(\theta_m|\hat{\theta}_{-m})| \leq \frac{C_7 np}{2} |\theta_m - \theta_m^*|^3 \leq \frac{C_7 np}{2} \left( \frac{|f^{(m)}(\theta_m^*|\hat{\theta}_{-m})|}{c_1 np} \right)^3.$$

By Pythagorean identity,  $\bar{\ell}_n^{(m)}(\hat{\theta}_m|\hat{\theta}_{-m}) = \bar{\ell}_n^{(m)}(\bar{\theta}_m|\hat{\theta}_{-m}) + \frac{1}{2}g^{(m)}(\theta_m^*|\hat{\theta}_{-m})(\hat{\theta}_m - \bar{\theta}_m)^2$ . Then

$$\begin{aligned} &\frac{1}{2}g^{(m)}(\theta_m^*|\hat{\theta}_{-m})(\hat{\theta}_m - \bar{\theta}_m)^2 \\ &= \bar{\ell}_n^{(m)}(\hat{\theta}_m|\hat{\theta}_{-m}) - \bar{\ell}_n^{(m)}(\bar{\theta}_m|\hat{\theta}_{-m}) \end{aligned}$$

$$\begin{aligned}
&\leq \ell_n^{(m)}(\widehat{\theta}_m|\widehat{\theta}_{-m}) - \bar{\ell}_n^{(m)}(\bar{\theta}_m|\widehat{\theta}_{-m}) + \frac{C_7 np}{2} \left( \frac{|f^{(m)}(\theta_m^*|\widehat{\theta}_{-m})|}{c_1 np} \right)^3 \\
&\leq \ell_n^{(m)}(\bar{\theta}_m|\widehat{\theta}_{-m}) - \bar{\ell}_n^{(m)}(\bar{\theta}_m|\widehat{\theta}_{-m}) + \frac{C_7 np}{2} \left( \frac{|f^{(m)}(\theta_m^*|\widehat{\theta}_{-m})|}{c_1 np} \right)^3 \\
&\leq 2 \frac{C_7 np}{2} \left( \frac{|f^{(m)}(\theta_m^*|\widehat{\theta}_{-m})|}{c_1 np} \right)^3.
\end{aligned}$$

Since  $g^{(m)}(\theta_m^*|\widehat{\theta}_{-m}) \geq c_1 np$ , we obtain the bound

$$(\widehat{\theta}_m - \bar{\theta}_m)^2 \leq \frac{2C_7}{c_1^4} \left( \frac{|f^{(m)}(\theta_m^*|\widehat{\theta}_{-m})|}{np} \right)^3.$$

Since  $|f^{(m)}(\theta_m^*|\widehat{\theta}_{-m})|$  has been shown to be bounded by the sum of (57)–(61), we have

$$(63) \quad |\widehat{\theta}_m - \bar{\theta}_m| \leq \delta \bar{\Delta}_m,$$

for some  $\delta = o(1)$  with probability at least  $1 - O(n^{-7}) - \exp(-\bar{\Delta}_m^{3/2} L np) - \exp(-\bar{\Delta}_m^2 np \times L \frac{np}{\log n})$  under the condition that  $\bar{\Delta}_m = o(1)$  and  $\frac{np}{\log n} \rightarrow \infty$ .

In the third step, we need to show that  $\frac{f^{(m)}(\theta_m^*|\widehat{\theta}_{-m})}{g^{(m)}(\theta_m^*|\theta_{-m}^*)}$  in the definition of  $\bar{\theta}_m$  can be replaced by  $\frac{f^{(m)}(\theta_m^*|\theta_{-m}^*)}{g^{(m)}(\theta_m^*|\theta_{-m}^*)}$  with a negligible error. By triangle inequality, we can bound  $|f^{(m)}(\theta_m^*|\widehat{\theta}_{-m}) - f^{(m)}(\theta_m^*|\theta_{-m}^*)|$  by the sum of (58), (59) and (61). Given that  $g^{(m)}(\theta_m^*|\theta_{-m}^*) \gtrsim np$ , we have

$$(64) \quad \frac{|f^{(m)}(\theta_m^*|\widehat{\theta}_{-m}) - f^{(m)}(\theta_m^*|\theta_{-m}^*)|}{g^{(m)}(\theta_m^*|\theta_{-m}^*)} \leq \delta \bar{\Delta}_m,$$

for some  $\delta = o(1)$  with probability at least  $1 - O(n^{-7}) - \exp(-\bar{\Delta}_m^2 np L \frac{np}{\log n})$  under the assumption that  $np L \bar{\Delta}_m^2 \rightarrow \infty$  and  $\frac{np}{\log n} \rightarrow \infty$ . Note that we can choose the same  $\delta$  to accommodate the two bounds (63) and (64). We also need to give a sharp approximation to  $g^{(m)}(\theta_m^*|\widehat{\theta}_{-m})$ . We have

$$\begin{aligned}
&|g^{(m)}(\theta_m^*|\widehat{\theta}_{-m}) - g^{(m)}(\theta_m^*|\theta_{-m}^*)| \\
&\leq |g^{(m)}(\theta_m^*|\widehat{\theta}_{-m}) - g^{(m)}(\theta_m^*|\theta_{-m}^{(m)} - a_m \mathbb{1}_{n-1})| \\
&\quad + |g^{(m)}(\theta_m^*|\theta_{-m}^{(m)} - a_m \mathbb{1}_{n-1}) - g^{(m)}(\theta_m^*|\theta_{-m}^*)| \\
&\leq \sqrt{\sum_{i \in [n] \setminus \{m\}} A_{mi}} \|\theta_{-m}^{(m)} - a_m \mathbb{1}_{n-1} - \widehat{\theta}_{-m}\| + p\sqrt{n} \|\theta_{-m}^{(m)} - a_m \mathbb{1}_{n-1} - \theta^*\| \\
&\quad + \sum_{i \in [n] \setminus \{m\}} (A_{mi} - p) |\theta_i^{(m)} - a_m - \theta_i^*|.
\end{aligned}$$

By Lemma 8.1, Lemma 8.6 and Lemma 8.7, the first two terms can be bounded by  $C_8 \sqrt{\frac{np}{L}}$  with probability at least  $1 - O(n^{-7})$ . To bound the third term, we can use Lemma 8.2, and then  $\sum_{i \in [n] \setminus \{m\}} (A_{mi} - p) |\theta_i^{(m)} - a_m - \theta_i^*|$  can be bounded by

$$C_8 \sqrt{p \log n} \|\theta_{-m}^{(m)} - a_m \mathbb{1}_{n-1} - \theta^*\| + C_8 \log n \|\theta_{-m}^{(m)} - a_m \mathbb{1}_{n-1} - \theta^*\|_\infty,$$

with probability at least  $1 - O(n^{-10})$ . By Lemma 8.6 and Lemma 8.7, the above display is at most  $C_9 \sqrt{\frac{\log n}{L}} + C_9 \frac{(\log n)^{3/2}}{\sqrt{npL}}$  with probability at least  $1 - O(n^{-7})$ . Combining our bounds,

we obtain

$$(65) \quad |g^{(m)}(\theta_m^*|\widehat{\theta}_{-m}) - g^{(m)}(\theta_m^*|\theta_{-m}^*)| \lesssim \sqrt{\frac{np}{L}} + \frac{(\log n)^{3/2}}{\sqrt{npL}}.$$

Since  $g^{(m)}(\theta_m^*|\theta_{-m}^*) \gtrsim np$ , we have

$$(66) \quad \frac{|g^{(m)}(\theta_m^*|\widehat{\theta}_{-m}) - g^{(m)}(\theta_m^*|\theta_{-m}^*)|}{g^{(m)}(\theta_m^*|\theta_{-m}^*)} \leq \delta,$$

for some  $\delta = o(1)$  with probability at least  $1 - O(n^{-7})$ . Note that we can choose the same  $\delta$  to accommodate the three bounds (63), (64) and (66).

In the last step, we will apply Lemma 3.1 with  $t = (1 - \eta)\theta_k^* + \eta\theta_{k+1}^*$  to complete the proof. Recall the definition of  $\eta$  in (54). For any  $i \leq k$ , we have

$$(67) \quad \begin{aligned} & \mathbb{P}(\widehat{\theta}_i \leq (1 - \eta)\theta_k^* + \eta\theta_{k+1}^*) \\ & \leq \mathbb{P}(\widehat{\theta}_i - \theta_i^* \leq -\eta(\theta_k^* - \theta_{k+1}^*) - (\theta_i^* - \theta_k^*)) \\ & \leq \mathbb{P}(\bar{\theta}_i - \theta_i^* \leq -(1 - \delta)\bar{\Delta}_i) + \mathbb{P}(|\bar{\theta}_i - \widehat{\theta}_i| > \delta\bar{\Delta}_i) \\ & \leq \mathbb{P}\left(-\frac{f^{(i)}(\theta_i^*|\theta_{-i}^*)}{g^{(i)}(\theta_i^*|\theta_{-i}^*)} \leq -(1 + \delta^2 - 3\delta)\bar{\Delta}_i\right) \\ & \quad + \mathbb{P}(|\bar{\theta}_i - \widehat{\theta}_i| > \delta\bar{\Delta}_i) \\ & \quad + \mathbb{P}\left(\frac{|g^{(i)}(\theta_i^*|\widehat{\theta}_{-i}) - g^{(i)}(\theta_i^*|\theta_{-i}^*)|}{g^{(i)}(\theta_i^*|\theta_{-i}^*)} > \delta\right) \\ & \quad + \mathbb{P}\left(\frac{|f^{(i)}(\theta_i^*|\widehat{\theta}_{-i}) - f^{(i)}(\theta_i^*|\theta_{-i}^*)|}{g^{(i)}(\theta_i^*|\theta_{-i}^*)} > \delta\bar{\Delta}_i\right) \\ & \leq \mathbb{P}\left(-\frac{f^{(i)}(\theta_i^*|\theta_{-i}^*)}{g^{(i)}(\theta_i^*|\theta_{-i}^*)} \leq -(1 - 3\delta)\bar{\Delta}_i\right) + O(n^{-7}) \\ & \quad + \exp(-\bar{\Delta}_i^{3/2}Lnp) + \exp\left(-\bar{\Delta}_i^2npL\frac{np}{\log n}\right), \end{aligned}$$

where the last inequality is due to (63), (64) and (66). Define the event

$$\mathcal{A}_i = \left\{ A : \left| \frac{\sum_{j \in [n] \setminus \{i\}} A_{ij} \psi(\theta_i^* - \theta_j^*) \psi(\theta_j^* - \theta_i^*)}{p \sum_{j \in [n] \setminus \{i\}} \psi(\theta_i^* - \theta_j^*) \psi(\theta_j^* - \theta_i^*)} - 1 \right| \leq \delta \right\}.$$

By Bernstein's inequality, we have  $\mathbb{P}(A \in \mathcal{A}_i^c) \leq O(n^{-7})$  for some  $\delta = o(1)$ . Again, we shall adjust the value of  $\delta$  so that (63), (64) and (66) are still true. We then have

$$(68) \quad \begin{aligned} & \mathbb{P}\left(-\frac{f^{(i)}(\theta_i^*|\theta_{-i}^*)}{g^{(i)}(\theta_i^*|\theta_{-i}^*)} \leq -(1 - 3\delta)\bar{\Delta}_i\right) \\ & \leq \sup_{A \in \mathcal{A}_i} \mathbb{P}\left(\frac{\sum_{j \in [n] \setminus \{i\}} A_{ij} (\bar{y}_{ij} - \psi(\theta_i^* - \theta_j^*))}{\sum_{j \in [n] \setminus \{i\}} A_{ij} \psi(\theta_i^* - \theta_j^*) \psi(\theta_j^* - \theta_i^*)} \leq -(1 - 3\delta)\bar{\Delta}_i | A\right) \\ & \quad + \mathbb{P}(A \in \mathcal{A}_i^c) \\ & \leq \sup_{A \in \mathcal{A}_i} \exp\left(-\frac{\frac{1}{2}(1 - 3\delta)^2 \bar{\Delta}_i^2 (L \sum_{j \in [n] \setminus \{i\}} A_{ij} \psi'(\theta_i^* - \theta_j^*))^2}{L \sum_{j \in [n] \setminus \{i\}} A_{ij} \psi'(\theta_i^* - \theta_j^*) + \frac{1-3\delta}{3} \bar{\Delta}_i L \sum_{j \in [n] \setminus \{i\}} A_{ij} \psi'(\theta_i^* - \theta_j^*)}\right) \end{aligned}$$

$$\begin{aligned}
& + O(n^{-7}) \\
(69) \quad & = \exp\left(-\frac{1+o(1)}{2}\bar{\Delta}_i^2 Lp \sum_{j \in [n] \setminus \{i\}} \psi'(\theta_i^* - \theta_j^*)\right) + O(n^{-7}) \\
(70) \quad & \leq \exp\left(-\frac{1+o(1)}{2}(\eta(\theta_k^* - \theta_{k+1}^*) + (\theta_i^* - \theta_k^*))^2 Lp \sum_{j \in [n] \setminus \{i\}} \psi'(\theta_i^* - \theta_j^*)\right) \\
& + O(n^{-7}) \\
(71) \quad & \leq \exp\left(-\frac{1+o(1)}{2}(\bar{\Delta} + (\theta_i^* - \theta_k^*))^2 Lp \sum_{j \in [n] \setminus \{i\}} \psi'(\theta_i^* - \theta_j^*)\right) + O(n^{-7}).
\end{aligned}$$

The bound (68) is by Bernstein's inequality. We then use the definition of  $\mathcal{A}_i$  to obtain the expression (69). To see why (70) is true, note that when  $\bar{\Delta}_i^2 = \sqrt{\frac{\log n}{np}}$ , the first term of (69) can be absorbed into  $O(n^{-7})$ . Finally, in (71), we have used the notation  $\bar{\Delta} = \min(\eta(\theta_k^* - \theta_{k+1}^*), (\frac{\log n}{np})^{1/4})$ . For each  $j \in [n]$ , define

$$h_j(t) = (\bar{\Delta} + t)^2 \psi'(t + \theta_k^* - \theta_j^*), \quad \text{for all } t \geq 0.$$

The derivative of this function is

$$h'_j(t) = (\bar{\Delta} + t) \psi'(t + \theta_k^* - \theta_j^*) [2 + (\bar{\Delta} + t)(1 - 2\psi(t + \theta_k^* - \theta_j^*))].$$

Since  $\max_{j,k} |\theta_k^* - \theta_j^*| = O(1)$ , we can find a sufficiently small constant  $c_2 > 0$ , such that  $h_j(t)$  is increasing on  $[0, c_2]$ . Moreover, there exists another small constant  $c_3 > 0$  such that  $\min_{t \in (c_2, \kappa]} h_j(t) \geq c_3$ . With this fact, we can bound the exponent of (71) as

$$\begin{aligned}
& (\bar{\Delta} + (\theta_i^* - \theta_k^*))^2 Lp \sum_{j \in [n] \setminus \{i\}} \psi'(\theta_i^* - \theta_j^*) \\
& \geq Lp \sum_{j \in [n] \setminus \{i\}} \min(\bar{\Delta}^2 \psi'(\theta_k^* - \theta_j^*), c_3) \\
(72) \quad & \geq Lp(k-1) \min(\bar{\Delta}^2 \psi'(\theta_1^* - \theta_k^*), c_3) + Lp(n-k) \min(\bar{\Delta}^2 \psi'(\theta_k^* - \theta_n^*), c_3) \\
& = Lp \bar{\Delta}^2 ((k-1) \psi'(\theta_1^* - \theta_k^*) + (n-k) \psi'(\theta_k^* - \theta_n^*)) \\
& \geq (1+o(1)) Lp \min\left(\eta^2 \Delta^2, \sqrt{\frac{\log n}{np}}\right) \frac{n}{V(\kappa)}
\end{aligned}$$

where the equality (72) uses the fact that  $\bar{\Delta} \rightarrow 0$ . Therefore, we can further bound (71) as

$$\begin{aligned}
& \exp\left(-\frac{1+o(1)}{2} Lp \min\left(\eta^2 \Delta^2, \sqrt{\frac{\log n}{np}}\right) \frac{n}{V(\kappa)}\right) + O(n^{-7}) \\
& \leq \exp\left(-\frac{(1+o(1)) \eta^2 \Delta^2 npL}{2V(\kappa)}\right) + O(n^{-7}).
\end{aligned}$$

The last inequality holds because when  $\min(\eta^2 \Delta^2, \sqrt{\frac{\log n}{np}}) = \sqrt{\frac{\log n}{np}}$ , the first term becomes  $\exp(-\frac{(1+o(1))L\sqrt{np\log n}}{2V(\kappa)})$ , which can be absorbed by  $O(n^{-7})$ . Since  $\exp(-\bar{\Delta}_i^{3/2} Lnp) + \exp(-\bar{\Delta}_i^2 npL \frac{np}{\log n}) \leq \exp(-\frac{(1+o(1))\eta^2 \Delta^2 npL}{2V(\kappa)}) + O(n^{-7})$ , we have

$$(73) \quad \mathbb{P}(\hat{\theta}_i \leq (1-\eta)\theta_k^* + \eta\theta_{k+1}^*) \leq \exp\left(-\frac{(1-\delta')\eta^2 \Delta^2 npL}{2V(\kappa)}\right) + O(n^{-7}),$$

with some  $\delta' = o(1)$  for all  $i \leq k$ . With a similar argument, we also have

$$(74) \quad \mathbb{P}(\hat{\theta}_i \geq (1 - \eta)\theta_k^* + \eta\theta_{k+1}^*) \leq \exp\left(-\frac{(1 - \delta')(1 - \eta)^2 \Delta^2 npL}{2V(\kappa)}\right) + O(n^{-7}),$$

for all  $i \geq k + 1$ . It can be checked that the  $\delta'$  above is independent of the  $\bar{\delta}$  in the definition of  $\eta$ . Now we can choose  $\eta$  as in (54) with  $\bar{\delta} = \delta'$ . By Lemma 3.1, we have

$$\begin{aligned} \mathbb{E}H_k(\hat{r}, r^*) &\leq \exp\left(-\frac{(1 - \bar{\delta})\eta^2 \Delta^2 npL}{2V(\kappa)}\right) \\ &\quad + \frac{n - k}{k} \exp\left(-\frac{(1 - \bar{\delta})(1 - \eta)^2 \Delta^2 npL}{2V(\kappa)}\right) + O(n^{-7}) \\ &\leq 2 \exp\left(-\frac{1}{2} \left( \frac{\sqrt{(1 - \bar{\delta})\text{SNR}}}{2} - \frac{1}{\sqrt{(1 - \bar{\delta})\text{SNR}}} \log \frac{n - k}{k} \right)^2\right) + O(n^{-7}). \end{aligned}$$

By Markov's inequality, the above bound implies

$$H_k(\hat{r}, r^*) \leq \exp\left(-\frac{1}{2} \left( \frac{\sqrt{(1 - \delta_1)\text{SNR}}}{2} - \frac{1}{\sqrt{(1 - \delta_1)\text{SNR}}} \log \frac{n - k}{k} \right)^2\right) + O(n^{-6}),$$

for some  $\delta_1 = o(1)$  with high probability. One can take, for example,

$$\delta_1 = \bar{\delta} + \frac{1}{\frac{\sqrt{(1 - \bar{\delta})\text{SNR}}}{2} - \frac{1}{\sqrt{(1 - \bar{\delta})\text{SNR}}} \log \frac{n - k}{k}}.$$

When  $O(n^{-6})$  dominates the bound, we have  $H_k(\hat{r}, r^*) = O(n^{-6})$ , which implies  $H_k(\hat{r}, r^*) = 0$  since  $H_k(\hat{r}, r^*) \in \{0, (2k)^{-1}, 2(2k)^{-1}, 3(2k)^{-1}, \dots, 1\}$ . Therefore, we always have

$$H_k(\hat{r}, r^*) \leq 2 \exp\left(-\frac{1}{2} \left( \frac{\sqrt{(1 - \delta_1)\text{SNR}}}{2} - \frac{1}{\sqrt{(1 - \delta_1)\text{SNR}}} \log \frac{n - k}{k} \right)^2\right),$$

with high probability for some  $\delta_1 = o(1)$ . The proof is complete.  $\square$

**PROOF OF THEOREM 3.3.** With some rearrangements, the condition is equivalent to

$$\frac{npL\Delta^2}{2(1 + \epsilon)V(\kappa)} \left( \frac{1}{2} - \frac{(1 + \epsilon)V(\kappa)}{npL\Delta^2} \log \frac{n - k}{k} \right)^2 > \log k.$$

Since  $\epsilon$  is a constant, it implies

$$\frac{npL\Delta^2}{2V(\kappa)} \left( \frac{1}{2} - \frac{V(\kappa)}{(1 - \delta)npL\Delta^2} \log \frac{n - k}{k} \right)^2 > (1 + \epsilon) \log k,$$

for any  $\delta = o(1)$ . Therefore,  $H_k(\hat{r}, r^*) = o(k^{-1})$  when  $k \rightarrow \infty$ . Given the fact that  $H_k(\hat{r}, r^*) \in \{0, (2k)^{-1}, 2(2k)^{-1}, 3(2k)^{-1}, \dots, 1\}$ , we must have  $H_k(\hat{r}, r^*) = 0$ . When  $k = O(1)$ , the condition implies  $\frac{npL\Delta^2}{2V(\kappa)} > (1 + \epsilon') \log n$  for some constant  $\epsilon' > 0$ . This leads to the fact that  $(\frac{1}{2} - \frac{V(\kappa)}{(1 - \delta)npL\Delta^2} \log \frac{n - k}{k})^2 > c_1$  for some constant  $c_1 > 0$ . Therefore,  $H_k(\hat{r}, r^*) = o(1) = o(k^{-1})$ , which implies  $H_k(\hat{r}, r^*) = 0$ .  $\square$

**Funding.** The first and second authors were supported in part by NSF CAREER award DMS-1847590 and NSF Grant CCF-1934931.

The third author was supported in part by NSF Grant DMS-2112988.

## SUPPLEMENTARY MATERIAL

**Supplement to “Partial recovery for top- $k$  ranking: Optimality of MLE and suboptimality of the spectral method”** (DOI: [10.1214/21-AOS2166SUPP](https://doi.org/10.1214/21-AOS2166SUPP); .pdf). The supplement [5] includes all the technical proofs. In Appendix A, we first give proofs for all the results established in Section 4: Theorem 4.1, Theorem 4.2 and Theorem 4.3. After that, we prove Theorem 3.4 and Theorem 6.1 in Appendix B. We then include the proofs of Theorem 7.1 and Theorem 7.2 in Appendix C, the proof of Lemma 8.5 in Appendix D and the proofs of all the other technical lemmas in Appendix E. The count method is discussed in Appendix F.

## REFERENCES

- [1] ABBE, E., FAN, J., WANG, K. and ZHONG, Y. (2020). Entrywise eigenvector analysis of random matrices with low expected rank. *Ann. Statist.* **48** 1452–1474. [MR4124330](https://doi.org/10.1214/19-AOS1854) <https://doi.org/10.1214/19-AOS1854>
- [2] BORDA, J. D. (1784). Mémoire sur les élections au scrutin. Histoire de l’Académie Royale des Sciences pour 1781 (Paris, 1784).
- [3] BRADLEY, R. A. and TERRY, M. E. (1952). Rank analysis of incomplete block designs. I. The method of paired comparisons. *Biometrika* **39** 324–345. [MR0070925](https://doi.org/10.2307/2334029) <https://doi.org/10.2307/2334029>
- [4] BUTUCEA, C., NDAOUD, M., STEPANOVA, N. A. and TSYBAKOV, A. B. (2018). Variable selection with Hamming loss. *Ann. Statist.* **46** 1837–1875. [MR3845003](https://doi.org/10.1214/17-AOS1572) <https://doi.org/10.1214/17-AOS1572>
- [5] CHEN, P., GAO, C. and ZHANG, A. Y. (2022). Supplement to “Partial Recovery for Top- $k$  Ranking: Optimality of MLE and SubOptimality of the Spectral Method.” <https://doi.org/10.1214/21-AOS2166SUPP>
- [6] CHEN, X., GOPI, S., MAO, J. and SCHNEIDER, J. (2017). Competitive analysis of the top- $K$  ranking problem. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms* 1245–1264. SIAM, Philadelphia, PA. [MR3627810](https://doi.org/10.1137/1.9781611974782.81) <https://doi.org/10.1137/1.9781611974782.81>
- [7] CHEN, Y., FAN, J., MA, C. and WANG, K. (2019). Spectral method and regularized MLE are both optimal for top- $K$  ranking. *Ann. Statist.* **47** 2204–2235. [MR3953449](https://doi.org/10.1214/18-AOS1745) <https://doi.org/10.1214/18-AOS1745>
- [8] CHEN, Y. and SUH, C. (2015). Spectral mle: Top- $k$  rank aggregation from pairwise comparisons. In *International Conference on Machine Learning* 371–380.
- [9] CHERNOZHUKOV, V., CHETVERIKOV, D. and KATO, K. (2013). Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors. *Ann. Statist.* **41** 2786–2819. [MR3161448](https://doi.org/10.1214/13-AOS1161) <https://doi.org/10.1214/13-AOS1161>
- [10] COPELAND, A. (1951). A ‘reasonable’ social welfare function, Seminar on mathematics in social sciences, University of Michigan. Cited Indirectly from Its Mention by Luce and Raiffa (1957) 358.
- [11] COSSOCK, D. and ZHANG, T. (2006). Subset ranking using regression. In *Learning Theory. Lecture Notes in Computer Science* **4005** 605–619. Springer, Berlin. [MR2280634](https://doi.org/10.1007/11776420_44) [https://doi.org/10.1007/11776420\\_44](https://doi.org/10.1007/11776420_44)
- [12] DWORK, C., KUMAR, R., NAOR, M. and SIVAKUMAR, D. (2001). Rank aggregation methods for the web. In *Proceedings of the 10th International Conference on World Wide Web* 613–622.
- [13] ERDŐS, P. and RÉNYI, A. (1960). On the evolution of random graphs. *Magy. Tud. Akad. Mat. Kut. Intéz. Közl.* **5** 17–61. [MR0125031](https://doi.org/10.1007/11776420_44)
- [14] JANG, M., KIM, S., SUH, C. and OH, S. (2016). Top- $k$  ranking from pairwise comparisons: When spectral ranking is optimal. ArXiv preprint. Available at [arXiv:1603.04153](https://arxiv.org/abs/1603.04153).
- [15] JANG, M., KIM, S., SUH, C. and OH, S. (2017). Optimal sample complexity of  $m$ -wise data for top- $k$  ranking. In *Advances in Neural Information Processing Systems* 1686–1696.
- [16] LÖFFLER, M., ZHANG, A. Y. and ZHOU, H. H. (2021). Optimality of spectral clustering in the Gaussian mixture model. *Ann. Statist.* **49** 2506–2530. [MR4338373](https://doi.org/10.1214/20-aos2044) <https://doi.org/10.1214/20-aos2044>
- [17] LU, Y. and ZHOU, H. H. (2016). Statistical and computational guarantees of Lloyd’s algorithm and its variants. ArXiv preprint. Available at [arXiv:1612.02099](https://arxiv.org/abs/1612.02099).
- [18] LUCE, R. D. (2012). *Individual Choice Behavior: A Theoretical Analysis*. Courier Corporation.
- [19] MOTEGI, S. and MASUDA, N. (2012). A network-based dynamical ranking system for competitive sports. *Sci. Rep.* **2** 904. <https://doi.org/10.1038/srep00904>
- [20] NDAOUD, M. and TSYBAKOV, A. B. (2020). Optimal variable selection and adaptive noisy compressed sensing. *IEEE Trans. Inf. Theory* **66** 2517–2532. [MR4087700](https://doi.org/10.1109/TIT.2020.2965738) <https://doi.org/10.1109/TIT.2020.2965738>
- [21] NEGAHBAN, S., OH, S. and SHAH, D. (2017). Rank centrality: Ranking from pairwise comparisons. *Oper. Res.* **65** 266–287. [MR3613103](https://doi.org/10.1287/opre.2016.1534) <https://doi.org/10.1287/opre.2016.1534>

- [22] SHA, L., LUCEY, P., YUE, Y., CARR, P., ROHLF, C. and MATTHEWS, I. (2016). Chalkboarding: A new spatiotemporal query paradigm for sports play retrieval. In *Proceedings of the 21st International Conference on Intelligent User Interfaces* 336–347.
- [23] SHAH, N. B. and WAINWRIGHT, M. J. (2017). Simple, robust and optimal ranking from pairwise comparisons. *J. Mach. Learn. Res.* **18** 199. [MR3827087](#)
- [24] TROPP, J. A. (2015). An introduction to matrix concentration inequalities. *Found. Trends Mach. Learn.* **8** 1–230.
- [25] ZHANG, A. Y. and ZHOU, H. H. (2016). Minimax rates of community detection in stochastic block models. *Ann. Statist.* **44** 2252–2280. [MR3546450](#) <https://doi.org/10.1214/15-AOS1428>