

# Matrix Completion with Covariate Information and Informative Missingness

**Huaqing Jin\***

*Department of Statistics and Actuarial Science  
The University of Hong Kong, Hong Kong*

HUAQING.JIN@UCSF.EDU

**Yanyuan Ma**

*Department of Statistics  
Pennsylvania State University*

YZM63@PSU.EDU

**Fei Jiang<sup>†</sup>**

*Department of Epidemiology and Biostatistics  
The University of California, San Francisco*

FEI.JIANG@UCSF.EDU

**Editor:** Benjamin Marlin

## Abstract

We study the problem of matrix completion when the missingness of the matrix entries is dependent on the unobserved response values themselves and hence the missingness itself is informative. Furthermore, we allow to take into account the covariate information to establish its relation with the response and hence enable prediction. We devise a novel procedure to simultaneously complete the partially observed matrix and assess the covariate effect. Allowing the matrix dimensions as well as the number of covariates to grow ultra-high, under the classic low-rank matrix and sparse covariate effect assumptions, we rigorously establish the statistical guarantee of our procedure and the algorithmic convergence. The method is demonstrated via simulation studies and is used to analyze a Yelp data set and a MovieLens data set.

keywords: matrix completion, informative missing, low rank, sparse, tensor

## 1. Motivation and Introduction

Matrix completion has become a popular research topic in both statistics and computer science, mainly driven by the commercial need from online providers. Consider a concrete example from Yelp, where the scores from a total of  $n$  customers evaluating  $m$  restaurants are of interest. Let  $Y_{ij}$  be an indicator denoting whether the evaluation of the  $i$ th subject regarding the  $j$ th restaurant is positive, where  $i = 1, \dots, n, j = 1, \dots, m$ . Let  $\mathbf{Y}$  be the collection of  $Y_{ij}$ 's, i.e.  $\mathbf{Y}$  is a  $n \times m$  matrix. Obviously, not everyone will evaluate every restaurant, so many elements in  $\mathbf{Y}$  are missing. Let  $R_{ij}$  denote the corresponding missingness index and let  $\mathbf{R}$  be the corresponding matrix formed by the collection of  $R_{ij}$ 's. It is sensible to suspect that the very fact of missingness is related to the potential evalua-

---

\*. Currently in Department of Radiology and Biomedical Imaging, UCSF.

†. Corresponding Author.

tion score itself—in fact, some individuals may dislike some restaurants so much that they refuse to visit them hence will naturally not evaluate them. Thus, if they had visited and evaluated, the feedback would be likely negative. This naturally leads to the nonignorable missing mechanism, where  $Y_{ij}$  and  $R_{ij}$  are dependent on each other. To recover the missing entries enables us to predict the missing evaluating scores, but is challenging because it is impossible to infer the distribution of  $Y_{ij}$  based on the biased sample formed by the observed entries. The issue can be resolved when there exist additional baseline covariates for each observation, while these covariates are independent of  $R_{ij}$  given  $Y_{ij}$ . When such covariate information is available, we aim to predict the missing entries of  $\mathbf{Y}$  based on the relation between  $Y_{ij}$  and the covariates.

Luckily, in the Yelp data, a large number of covariates have been collected, which include but are not limited to the covariates of a restaurant such as size, price, open hours, and traits of a customer such as the number of his/her friends and the number of his/her restaurant reviews. These covariates are fully observed, which makes the parameter estimation feasible even though  $Y_{ij}$  and  $R_{ij}$  are correlated. On the other hand, the covariates are of the ultra-high dimension, which brings difficulties in parameter estimation. To take into account both the high dimensional covariates and the nonignorable missingness nature of such data, we impose the familiar sparsity and low-rank constraints on the relation between  $Y_{ij}$  and the covariate vector  $\mathbf{X}_{ij}$ . Furthermore, in practice, the missing-data mechanism is often not well understood. Thus, we consider statistical methods that do not require specification of the mechanism. This leads us to the matrix completion problem with ultra-high dimensional covariate and nonignorable missingness.

In the remaining text, we describe our statistical procedure in Section 2 and establish the finite sample statistical properties in Section 3. We consider the computational issue both algorithmically and theoretically in Section 4. Simulated examples are demonstrated in Section 5 and we analyze both a Yelp data set and a MovieLens data set in Section 6. We conclude the paper in Section 8, and relegate the proof details to an Appendix.

## 2. Methodology

We now present the probability model used to construct likelihood for the parameter estimation. Let  $\mathbf{X}_{ij}$  be a  $p$ -dimensional covariate vector. Furthermore let the  $k$ th element of  $\mathbf{X}_{ij}$  be  $X_{ijk}$  and let  $\mathbf{X}$  be the  $n \times m \times p$  tensor whose  $(i, j, k)$ th element is  $X_{ijk}$ . Let  $\mathbf{Y}$  be the  $n \times m$  matrix whose  $(i, j)$ th element is  $Y_{ij}$ . We work in a very flexible model class, the generalized linear model, where the conditional density of  $Y_{ij}$  given  $\mathbf{X}_{ij}$  is  $f(Y_{ij}, \Theta_{0ij} + \beta_0^\top \mathbf{X}_{ij})$ , where  $\beta_0$  is the true effect from  $\mathbf{X}_{ij}$  on  $Y_{ij}$ ,  $\Theta_{0ij}$  is the  $(i, j)$ th element of  $\Theta_0$ . Here  $\Theta_0$  is the true intercept matrix which is a low rank  $n \times m$  matrix with rank  $r$ . Because it is expected that only a small number of covariates may affect the outcome, we assume the widely adopted sparsity assumption on  $\beta_0$ , and assume the sparseness of  $\beta_0$  is  $s$ , i.e.  $\|\beta_0\|_0 = s$ . When we focus on the estimation of the parameter  $\beta_0$ , we also assume  $\sum_k \beta_{0k} \mathbf{X}^k$  to be sparse to ensure identifiability, where  $\mathbf{X}^k$  is a  $n \times m$  matrix whose  $(i, j)$ th element is  $X_{ijk}$ . Of course, when we only aim at prediction, then we do not make this assumption. Here despite of the sparsity and the low-rank assumptions, we do not restrict  $s, r$  to be finite. In other words, we allow  $s, r$  to grow with the sample size, although at a slower rate than  $p$  and  $n, m$ . For notational simplicity, we denote  $d = \sqrt{mn}$ .

Under the informative missing setting,  $Y_{ij}$  and  $R_{ij}$  are dependent. This implies that the expectation of the logarithm of the observed data likelihood  $\sum_{i=1}^n \sum_{j=1}^m R_{ij} \log f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})$  is no longer maximized at the true values  $\beta_0, \Theta_0$ , because  $E\{R_{ij} \log f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})\}$  does not equal to  $E(R_{ij} | \mathbf{X}_{ij}) E\{\log f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})\}$ , which would have been the case for the noninformative missing. Hence when estimating  $\Theta_{ij}$  and  $\beta$ , the information in  $R_{ij}$  cannot be ignored, and the loss function constructed from this likelihood cannot be used to estimate the parameters. On the other hand, because  $\mathbf{X}_{ij}$  is fully observed, under the assumption that  $\mathbf{X}_{ij}$  is independent of  $R_{ij}$  when  $Y_{ij}$  is given, even though  $Y_{ij}$  may not be available, we still can use the conditional distribution of  $\text{pr}(\mathbf{X}_{ij} | Y_{ij}, \Theta_0, \beta_0)$  to construct a pseudo-likelihood. In fact, the conditional independence between  $\mathbf{X}_{ij}$  and  $R_{ij}$  given  $Y_{ij}$  is not essential and can be relaxed. As long as part of the covariates in  $\mathbf{X}_{ij}$  are independent of  $R_{ij}$  given  $Y_{ij}$  and the remaining part in  $\mathbf{X}_{ij}$ , we can still construct a pseudo-likelihood. The essential benefit of considering a pseudo-likelihood is in eliminating the sampling bias. Intuitively, the observed portion of the data, i.e. the complete data, form a biased sample of the hypothetical full data, caused by the sampling bias in  $Y_{ij}$ 's. Through conditioning on  $Y_{ij}$ 's, the problem on the surface becomes to study the dependence of  $\mathbf{X}_{ij}$  given  $Y_{ij}$ , which does not relate to how  $Y_{ij}$  is sampled any more since the sampling of  $Y_{ij}$  now becomes a design issue hence could be done in any way we like. To help better understanding this issue, consider a standard regression problem where the response variable is named  $\mathbf{X}$  and the covariates named  $Y$ . In estimating the parameter involved in the model of  $\mathbf{X} | Y$ , we can construct various estimators based on the pairs of data  $(Y, \mathbf{X})$ 's without worrying how the covariates  $Y$ 's are sampled. In fact, even if the covariates  $Y$ 's are collected by design, say we only collected the covariates on a grid in a fixed region, it does not prevent us from obtaining a valid estimator based on the collected data. Here, in our context, we can view the non-missing  $Y_{ij}$ 's as the collected covariates, and the complete data  $(Y_{ij}, \mathbf{X}_{ij})$ 's as the observations in a standard regression problem to help grasp the intuition behind the pseudo likelihood method. We name this view point a covariate-dependent design scheme.

Specifically, starting from the conditional distribution of  $\mathbf{X}_{ij}$  given  $Y_{ij}$ , the average of the logarithm of the pseudo-likelihood, i.e. the conditional likelihood of the complete data, is written as

$$(mn)^{-1} \log \left[ \prod_{i=1}^n \prod_{j=1}^m \{\text{pr}(\mathbf{X}_{ij} | Y_{ij}, \Theta_0, \beta_0)\}^{R_{ij}} \right] \propto -(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m \mathcal{L}_{ij}(\Theta_0, \beta_0),$$

where  $\mathcal{L}_{ij}(\Theta, \beta) = R_{ij} \ell_{ij}(\Theta, \beta)$ ,

$$\ell_{ij}(\Theta, \beta) = - \left[ \log \{f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})\} - \log \left\{ \int f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X} \right\} \right],$$

and  $g(\cdot)$  is the joint Radon-Nikodym probability density function of  $\mathbf{X}_{ij}$ , i.e., the Radon-Nikodym derivative of the probability distribution of interest with respect to the dominating measure. We require that  $g(\cdot)$  is not a Dirac function. The dominating measure for the continuous component in  $\mathbf{X}_{ij}$  is the Lebesgue measure. The dominating measure for the discrete component in  $\mathbf{X}_{ij}$  is the counting measure, and therefore the integration with respect to the discrete component is the sum over its domain. Clearly, because  $\text{Pr}(R_{ij} = 1 | Y_{ij}, \mathbf{X}_{ij}) = \text{Pr}(R_{ij} = 1 | Y_{ij})$ ,  $E\{R_{ij} \ell_{ij}(\Theta, \beta) | Y_{ij}\} = E(R_{ij} | Y_{ij}) E\{\ell_{ij}(\Theta, \beta) | Y_{ij}\}$ . Therefore

the minimizer of  $E\{\ell_{ij}(\boldsymbol{\Theta}, \boldsymbol{\beta}) \mid Y_{ij}\}$  is the minimizer of  $E\{\mathcal{L}_{ij}(\boldsymbol{\Theta}, \boldsymbol{\beta}) \mid Y_{ij}\}$ . Furthermore,  $\Pr(R_{ij} = 1 \mid Y_{ij})$  and the marginal distribution of  $Y_{ij}$  do not contain information regarding the unknown parameters. Therefore taking expectation with respect to  $Y_{ij}$ , we minimize  $E\{\ell_{ij}(\boldsymbol{\Theta}, \boldsymbol{\beta})\}$  through minimizing  $E\{\mathcal{L}_{ij}(\boldsymbol{\Theta}, \boldsymbol{\beta})\}$ .

Now, let  $\mathcal{L}(\boldsymbol{\Theta}, \boldsymbol{\beta}) \equiv (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m \mathcal{L}_{ij}(\boldsymbol{\Theta}, \boldsymbol{\beta})$  be the negative pseudo-likelihood. Conditional on the observed  $Y_{ij}$ 's, the observed covariates form a conditional random sample even when the response missingness depends on the response value  $Y_{ij}$  itself. Furthermore, the expected logarithm of the pseudo-likelihood, i.e.  $E\{\mathcal{L}_{ij}(\boldsymbol{\Theta}, \boldsymbol{\beta})\}$ , is maximized at  $\boldsymbol{\Theta}_0$  and  $\boldsymbol{\beta}_0$ . This can be verified in two ways. The first way continues from the covariate design point of view. Because the pseudo-likelihood is the same as the conditional likelihood of  $\mathbf{X}_{ij} \mid Y_{ij}$  restricted to the complete data, while the complete data are obtained from a covariate-dependent (i.e.  $Y$  dependent) design scheme, hence the maximizer is the true regression parameters  $\boldsymbol{\Theta}_0$  and  $\boldsymbol{\beta}_0$  at infinite samples. The second way is through detailed mathematical computation. We can verify that the derivative of the log pseudo-likelihood has expectation zero at  $\boldsymbol{\Theta}_0$  and  $\boldsymbol{\beta}_0$ . Further, we can also verify that the log pseudo-likelihood is a convex function, so it is indeed maximized at  $\boldsymbol{\Theta}_0$  and  $\boldsymbol{\beta}_0$  at infinite samples. Hence, following standard M-estimation theory (Van Der Vaart and Wellner, 2000), maximizing the pseudo-likelihood will indeed lead to a valid statistical estimation procedure. Note that the function we maximize is not the same as the log-likelihood of the observed data, regardless we write the likelihood in terms of  $\mathbf{X}_{ij} \mid Y_{ij}$  or  $Y_{ij} \mid \mathbf{X}_{ij}$ , but is the conditional likelihood of the complete data, hence it is named pseudo-likelihood.

When the dimensions of  $\boldsymbol{\Theta}$  and  $\boldsymbol{\beta}$  are ultra-high, the value of the regression function can diverge to infinity if the entries of  $\boldsymbol{\Theta}$  and  $\boldsymbol{\beta}$  are not bounded. This is unreasonable when the response has finite mean. To control the magnitude of the regression function, we assume that  $\boldsymbol{\Theta}_0$  and  $\boldsymbol{\beta}_0$  have finite entries, and search the estimators in the feasible sets that  $\|\boldsymbol{\Theta}\|_{\max} \leq a$  and  $\|\boldsymbol{\beta}\|_{\infty} \leq a$  for a constant  $a > 0$ . Combining with the low-rank and sparse structures of  $\boldsymbol{\Theta}$  and  $\boldsymbol{\beta}$  respectively, we propose to estimate  $\boldsymbol{\Theta}_0, \boldsymbol{\beta}_0$ , through

$$(\hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\beta}}) = \underset{\|\boldsymbol{\Theta}\|_{\max} \leq a, \|\boldsymbol{\beta}\|_{\infty} \leq a}{\operatorname{argmin}} \mathcal{L}(\boldsymbol{\Theta}, \boldsymbol{\beta}) + \lambda_{\boldsymbol{\Theta}} \|\boldsymbol{\Theta}\|_* + \lambda_{\boldsymbol{\beta}} \|\boldsymbol{\beta}\|_1. \quad (1)$$

Here,  $\|\boldsymbol{\Theta}\|_*$  is the nuclear norm of the matrix  $\boldsymbol{\Theta}$ , which drives the resulting estimator towards a low-rank matrix. Furthermore,  $\|\boldsymbol{\beta}\|_1$  is the  $L_1$  norm of  $\boldsymbol{\beta}$ , which induces the sparseness of the estimated  $\boldsymbol{\beta}$ . In general, for any matrix  $\mathbf{A}$ , we use  $\|\mathbf{A}\|_F$ ,  $\|\mathbf{A}\|_{op}$ ,  $\|\mathbf{A}\|_*$ ,  $\|\mathbf{A}\|_1$ ,  $\|\mathbf{A}\|_{\infty}$  and  $\|\mathbf{A}\|_{\max}$  to denote the Frobenius norm, spectral norm, nuclear norm, matrix 1-norm, matrix sup-norm and the element-wise maximal norm of the matrix  $\mathbf{A}$ , respectively. We also use  $\|\mathbf{a}\|_2$ ,  $\|\mathbf{a}\|_{\infty}$ ,  $\|\mathbf{a}\|_1$  to denote the  $L_2$ ,  $L_{\infty}$  and  $L_1$  norm of the vector  $\mathbf{a}$ , respectively. We use boldface letters to denote vectors or matrices throughout the text.

### 3. Statistical Properties

#### 3.1 Additional Notation

We introduce some additional notations to facilitate the presentation of the theoretic properties. Some of their explicit forms are derived in Appendix A. Let  $f_2$  and  $f_{22}$  be the first and second derivative of  $f$  with respect to the second argument, respectively. Let  $\mathbf{S}(Y_{ij}, \mathbf{X}_{ij} \mid \boldsymbol{\Theta}, \boldsymbol{\beta})$  be the derivative of  $-\log\{f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X}_{ij})\}$  with respect to  $\mathbf{M} =$

Notation	Definition
$f_2$	The first derivative of $f$ with respect to the second argument.
$f_{22}$	The second derivative of $f$ with respect to the second argument.
$\mathbf{z}_{ij}$	An $n \times m$ matrix with its $(i, j)$ th entry 1 and all other entries 0.
$d_{H_2}$	Upper bound of $ H_2(y, \mathbf{x} \boldsymbol{\Theta}, \boldsymbol{\beta}) $ .
$d_{S_2}$	Upper bound of $ S_2(y, \mathbf{x} \boldsymbol{\Theta}, \boldsymbol{\beta}) $ .
$a$	Upper bound of $\ \boldsymbol{\beta}\ _\infty$ and $\ \boldsymbol{\Theta}\ _{\max}$ .
$c_0$	Upper bound of $\ \mathbf{X}_{ij}\ _1$ .
$\mathbf{X}^k$	$n \times m$ matrix whose $(i, j)$ th element is $X_{ijk}$ .
$\sigma_{1F}$	$32c_0^2a^2(d_{H_2} + d_{S_2}^2)$ .
$\sigma_{dF}$	$32a^2(d_{H_2} + d_{S_2}^2)$ .
$W_{ijk}(\boldsymbol{\Theta}, \boldsymbol{\beta})$	$\mathbf{e}_k^T \partial^2 \ell_{ij}(\boldsymbol{\Theta}, \boldsymbol{\beta}) / \partial \boldsymbol{\beta} \partial (\boldsymbol{\Theta}, \mathbf{z}_{ij})$ .
$d_{\mathbf{W}}$	Upper bound of $\sup_k \ \text{vec}\{\mathbf{R} \circ \mathbf{W}_k(\boldsymbol{\Theta}, \boldsymbol{\beta}_0)\}\ _1$ .
$d_{EX}$	Upper bound of $\sup_k  \sum_{i=1}^n \sum_{j=1}^m R_{ij} E\{S_2(Y_{ij}, \mathbf{X} \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) X_k   Y_{ij}\} $ .
$\mathbf{E}_{ij}, E_{ijk}$	$-R_{ij} \partial \ell_{ij}(\boldsymbol{\Theta}, \boldsymbol{\beta}) / \partial \boldsymbol{\beta}, \mathbf{e}_k^T \mathbf{E}_{ij}$ .
$\mathbf{S}_{\boldsymbol{\Theta}}(Y_{ij}, \mathbf{X}_{ij} \boldsymbol{\Theta}, \boldsymbol{\beta})$	Partial derivatives of $-\log\{f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X}_{ij})\}$ with respect to $\boldsymbol{\Theta}$ .
$\mathbf{S}_{\boldsymbol{\beta}}(Y_{ij}, \mathbf{X}_{ij} \boldsymbol{\Theta}, \boldsymbol{\beta})$	Partial derivatives of $-\log\{f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X}_{ij})\}$ with respect to $\boldsymbol{\beta}$ .
$\mathbf{S}(Y_{ij}, \mathbf{X}_{ij} \boldsymbol{\Theta}, \boldsymbol{\beta})$	$[\text{vec}\{\mathbf{S}_{\boldsymbol{\Theta}}(Y_{ij}, \mathbf{X}_{ij} \boldsymbol{\Theta}, \boldsymbol{\beta})\}^T, \mathbf{S}_{\boldsymbol{\beta}}(Y_{ij}, \mathbf{X}_{ij} \boldsymbol{\Theta}, \boldsymbol{\beta})^T]^T$ .
$S_2(Y_{ij}, \mathbf{X}_{ij} \boldsymbol{\Theta}, \boldsymbol{\beta})$	The first derivative of $-\log f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X}_{ij})$ with respect to $\boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X}_{ij}$ .
$H_2(Y_{ij}, \mathbf{X}_{ij} \boldsymbol{\Theta}, \boldsymbol{\beta})$	The second derivative of $-\log f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X}_{ij})$ with respect to $\boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X}_{ij}$ .
$\mathbf{H}(Y_{ij}, \mathbf{X}_{ij} \boldsymbol{\Theta}, \boldsymbol{\beta})$	The Hessian matrix of $-\log\{f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X}_{ij})\}$ with respect to $\{\text{vec}(\boldsymbol{\Theta})^T, \boldsymbol{\beta}^T\}^T$ .
$\mathbf{H}_{\boldsymbol{\Theta}}(Y_{ij}, \mathbf{X}_{ij} \boldsymbol{\Theta}, \boldsymbol{\beta})$	The second derivatives of $-\log f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X})$ with respect to $\text{vec}(\boldsymbol{\Theta})$ .
$\mathbf{H}_{\boldsymbol{\beta}}(Y_{ij}, \mathbf{X}_{ij} \boldsymbol{\Theta}, \boldsymbol{\beta})$	The second derivatives of $-\log f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X})$ with respect to $\text{vec}(\boldsymbol{\beta})$ .
$\mathbf{F}_{\boldsymbol{\beta}}(\mathbf{X}, \mathbf{Y} \boldsymbol{\Theta}, \boldsymbol{\beta})$	$\partial^2 \mathcal{L}(\boldsymbol{\Theta}, \boldsymbol{\beta}) / \partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T$ .
$\mathbf{F}_{\boldsymbol{\Theta}}(\mathbf{X}, \mathbf{Y} \boldsymbol{\Theta}, \boldsymbol{\beta})$	$\partial^2 \mathcal{L}(\boldsymbol{\Theta}, \boldsymbol{\beta}) / \partial \text{vec}(\boldsymbol{\Theta}) \partial \text{vec}(\boldsymbol{\Theta})^T$ .
$E\{\mathbf{F}_{\boldsymbol{\beta}}(\mathbf{X}, \mathbf{Y} \hat{\boldsymbol{\Theta}}, \boldsymbol{\beta}^*)\}$	$E\{\mathbf{F}_{\boldsymbol{\beta}}(\mathbf{X}, \mathbf{Y} \boldsymbol{\Theta}, \boldsymbol{\beta})\} _{\boldsymbol{\Theta}=\hat{\boldsymbol{\Theta}}, \boldsymbol{\beta}=\boldsymbol{\beta}^*}$ .
$E\{\mathbf{F}_{\boldsymbol{\Theta}}(\mathbf{X}, \mathbf{Y} \hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\beta}})\}$	$E\{\mathbf{F}_{\boldsymbol{\Theta}}(\mathbf{X}, \mathbf{Y} \boldsymbol{\Theta}, \boldsymbol{\beta})\} _{\boldsymbol{\Theta}=\hat{\boldsymbol{\Theta}}, \boldsymbol{\beta}=\hat{\boldsymbol{\beta}}}$ .

Table 1: Notation.

$\{\text{vec}(\boldsymbol{\Theta})^T, \boldsymbol{\beta}^T\}^T$ . Obviously  $\mathbf{S}(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta})$  is the score function. We further decompose the score function as  $\mathbf{S}(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta}) = [\text{vec}\{\mathbf{S}_{\boldsymbol{\Theta}}(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta})\}^T, \mathbf{S}_{\boldsymbol{\beta}}(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta})^T]^T$ , where  $\mathbf{S}_{\boldsymbol{\Theta}}(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta})$  and  $\mathbf{S}_{\boldsymbol{\beta}}(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta})$  are the partial derivatives of  $-\log\{f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X}_{ij})\}$  with respect to  $\boldsymbol{\Theta}$  and  $\boldsymbol{\beta}$  respectively. Let  $\mathbf{S}_{\boldsymbol{\Theta}, k, l}(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta}), \mathbf{S}_{\boldsymbol{\beta}, k}(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta})$  be the  $(k, l)$ th element and  $k$ th element of  $\mathbf{S}_{\boldsymbol{\Theta}}(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta})$  and  $\mathbf{S}_{\boldsymbol{\beta}}(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta})$ , respectively.

Further let  $S_2(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta}), H_2(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta})$  be the first and second derivative of the negative log-likelihood  $-\log f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X}_{ij})$  with respect to  $\boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X}_{ij}$  and let  $\mathbf{H}(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta})$  be the Hessian matrix of  $-\log\{f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X}_{ij})\}$  with respect to  $\{\text{vec}(\boldsymbol{\Theta})^T, \boldsymbol{\beta}^T\}^T$ . Let  $\mathbf{H}_{\boldsymbol{\Theta}}(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta})$  and  $\mathbf{H}_{\boldsymbol{\beta}}(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta})$  be the diagonal blocks of  $\mathbf{H}(Y_{ij}, \mathbf{X}_{ij}|\boldsymbol{\Theta}, \boldsymbol{\beta})$ , i.e. they are the second derivatives of  $-\log f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^T \mathbf{X})$  with respect to  $\text{vec}(\boldsymbol{\Theta})$  and  $\boldsymbol{\beta}$  respectively.

Let  $\mathbf{z}_{ij}$  be an  $n \times m$  matrix with its  $(i, j)$ th entry 1 and all other entries 0. Note that we can extract the  $(i, j)$ th entry of  $\boldsymbol{\Theta}$  using  $\mathbf{z}_{ij}$  and  $\boldsymbol{\Theta}$ . To simplify the notation, we define  $W_{ijk}(\boldsymbol{\Theta}, \boldsymbol{\beta}) \equiv \mathbf{e}_k^T \partial^2 \ell_{ij}(\boldsymbol{\Theta}, \boldsymbol{\beta}) / \partial \boldsymbol{\beta} \partial (\boldsymbol{\Theta}, \mathbf{z}_{ij})$ . Then  $\mathbf{e}_k^T \partial^2 \mathcal{L}_{ij}(\boldsymbol{\Theta}, \boldsymbol{\beta}) / \partial \boldsymbol{\beta} \partial (\boldsymbol{\Theta}, \mathbf{z}_{ij}) = R_{ij} W_{ijk}(\boldsymbol{\Theta}, \boldsymbol{\beta})$ . Let  $\mathbf{W}_k(\boldsymbol{\Theta}, \boldsymbol{\beta})$  be the  $n \times m$  matrix with its  $(i, j)$ th element  $W_{ijk}(\boldsymbol{\Theta}, \boldsymbol{\beta})$ . Let  $\mathbf{E}_{ij} \equiv -R_{ij} \partial \ell_{ij}(\boldsymbol{\Theta}, \boldsymbol{\beta}) / \partial \boldsymbol{\beta}$  and  $E_{ijk} = \mathbf{e}_k^T \mathbf{E}_{ij}$ . We also write  $\mathbf{X}^k$  as the  $n \times m$  matrix whose  $(i, j)$ th element is  $X_{ijk}$ , i.e.  $\mathbf{X}^k$  is the  $k$ th slice of the tensor  $\mathbf{X}$ . Define  $\mathbf{F}_{\boldsymbol{\beta}}(\mathbf{X}, \mathbf{Y}|\boldsymbol{\Theta}, \boldsymbol{\beta}) =$

$\partial^2 \mathcal{L}(\Theta, \beta) / \partial \beta \partial \beta^T$ , and  $\mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y} | \Theta, \beta) = \partial^2 \mathcal{L}(\Theta, \beta) / \partial \text{vec}(\Theta) \partial \text{vec}(\Theta)^T$ . Specifically,

$$\begin{aligned} & \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \Theta, \beta) \\ = & (mn)^{-1} \left( \sum_{i=1}^n \sum_{j=1}^m R_{ij} [H_2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \mathbf{X}_{ij} \mathbf{X}_{ij}^T - E\{H_2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \mathbf{X}_{ij} \mathbf{X}_{ij}^T | Y_{ij}\}] \right. \\ & \left. + E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \mathbf{X}_{ij} \mathbf{X}_{ij}^T | Y_{ij}\} - E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2} \right), \end{aligned}$$

and

$$\begin{aligned} & \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y} | \Theta, \beta) \\ = & (mn)^{-1} \left( \sum_{i=1}^n \sum_{j=1}^m R_{ij} [H_2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T \right. \\ & - E\{H_2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\}] \\ & + E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \\ & \left. - E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \text{vec}(\mathbf{z}_{ij}) | Y_{ij}\}^{\otimes 2} \right), \end{aligned}$$

where all the expectations are with respect to the true distribution of  $\mathbf{X}_{ij}$  given  $Y_{ij}$ .

We use  $\beta^*$  to denote a point on the line between  $\hat{\beta}$  and  $\beta_0$  and  $\Theta^*$  to denote a point between  $\hat{\Theta}$  and  $\Theta_0$ , respectively. We summarize the notations in Table 1.

### 3.2 Conditions

To obtain the statistical convergence properties of our estimator, we assume the following regularity conditions.

(C1)  $\|\Theta_0\|_{\max} \leq a$ ,  $\|\beta_0\|_{\infty} \leq a$  for a constant  $a > 0$ .

(C2) The matrix  $\sum_{k=1}^p \beta_{0k} \mathbf{X}^k$  is sparse and the matrix  $\Theta$  has low rank. Specifically,  $\inf_{\rho > 0} \omega_1(\rho) \omega_2(\rho) < 1$ . Here,

$$\begin{aligned} \omega_1(\rho) & \equiv \max \left\{ \rho \left\| \text{sign} \left( \sum_{k=1}^p \beta_{0k} \mathbf{X}^k \right) \right\|_{1 \rightarrow 1}, \rho^{-1} \left\| \text{sign} \left( \sum_{k=1}^p \beta_{0k} \mathbf{X}^k \right) \right\|_{\infty \rightarrow \infty} \right\}, \\ \omega_2(\rho) & \equiv \rho^{-1} \|\mathbf{U} \mathbf{U}^T\|_{\max} + \rho \|\mathbf{V} \mathbf{V}^T\|_{\max} + \|\mathbf{U}\|_{2 \rightarrow \infty} \|\mathbf{V}\|_{2 \rightarrow \infty}, \end{aligned}$$

where for any  $n \times m$  matrix  $\mathbf{M}$ ,

$$\|\mathbf{M}\|_{p \rightarrow q} \equiv \max (\|\mathbf{M} \mathbf{v}\|_q, \mathbf{v} \in \mathbb{R}^m, \|\mathbf{v}\|_p \leq 1),$$

and  $\mathbf{U}, \mathbf{V}$  are the left and right singular vectors of  $\Theta_0$ .

(C3)  $\sup_k \|\text{vec}\{\mathbf{R} \circ \mathbf{W}_k(\Theta, \beta_0)\}\|_1 \leq d_{\mathbf{W}}$  for  $\Theta$  that satisfies  $\|\Theta\|_{\max} \leq a$ . Here,  $d_{\mathbf{W}} > 0$  and  $a$  is the constant defined in Condition (C1).

(C4) Assume

$$\sup_k \left| \sum_{i=1}^n \sum_{j=1}^m R_{ij} E\{S_2(Y_{ij}, \mathbf{X} | \Theta_0, \beta_0) X_k | Y_{ij}\} \right| \leq d_{E\mathbf{X}}.$$

Here  $d_{E\mathbf{X}} > 0$ . Here  $d_{E\mathbf{X}}$  does not need to be bounded.

(C5) Assume  $f(Y_{ij}, \Theta_{0ij} + \beta_0^T \mathbf{X}_{ij})$  is a continuously differentiable function with respect to  $\Theta_{0ij}$  and  $\beta_0$ . Furthermore,  $S_2(Y_{ij}, \mathbf{X}_{ij} | \Theta_0, \beta_0)$  is a sub-Gaussian random variable.  $\|\mathbf{X}_{ij}\|_1 \leq c_0$  for  $c_0 > 0$ . This implies  $\|\mathbf{X}_{ij}\|_2 \leq c_0$ .

(C6)  $|H_2(y, \mathbf{x} | \Theta, \beta)|$  is bounded uniformly by  $d_{H_2}$  and  $|S_2(y, \mathbf{x} | \Theta, \beta)|$  is bounded uniformly by  $d_{S_2}$  for all  $y, \mathbf{x}$ , and  $\Theta, \beta$  in the feasible set that  $\|\Theta\|_{\max} \leq a, \|\beta\|_{\infty} \leq a$ . Here  $d_{H_2} > 0$  and  $d_{S_2} > 0$  satisfy  $c_0^2 a^2 (d_{H_2} + d_{S_2}^2) \sqrt{\log\{\max(p, mn)\}} / (mn) \rightarrow 0$  and  $a^2 (d_{H_2} + d_{S_2}^2) \sqrt{d \log(d) / (mn)} \rightarrow 0$ .

Condition (C1) bounds the supnorms of  $\Theta_0$  and  $\beta_0$  so that the true parameters fall in the feasible set defined in (1). Let  $\mathbf{E} = \sum_{k=1}^p \beta_{0k} \mathbf{X}^k$ , then Condition (C2) is a standard identifiability condition for the matrix completion problem to identify  $\Theta_0$  and  $\mathbf{E}$ , where the regression function is the sum of a low-rank matrix and a sparse matrix (Hsu et al., 2011). Here, we set  $\Theta$  to be low rank and allow  $\mathbf{E} = \sum_{k=1}^p \beta_{0k} \mathbf{X}^k$  to be sparse because the covariate matrices in the social media data are often naturally sparse, where the majority of them are dictionary variables with a small number of nonzero entries (Robin et al., 2018). Without this condition, there will not be a clear definition of  $\beta_0$ , hence it is impossible to derive the distance  $\|\hat{\beta} - \beta_0\|_2$ . Condition (C2) needs to hold only when we are interested in estimating  $\Theta_0$  and  $\beta_0$ . If the goal is to predict the missing entries in  $\mathbf{Y}$ , Condition (C2) is not necessary. Note that we also assume sparsity on  $\beta_0$ , which is customary in treating high dimensional regression and does not pertain to the matrix completion problem alone. Condition (C3) ensures that the objective function has bounded second derivative with respect to  $\beta$  in  $L_1$  norm. Conditions (C4), (C5) and (C6) provide the upper bounds of the high dimensional covariate, score functions and Hessian matrix, respectively. It is important to note that  $c_0, d_{H_2}$  and  $d_{S_2}$  are not necessarily bounded. Their growing rates jointly determine the convergence of  $\hat{\beta}$  and  $\hat{\Theta}$  as shown in Theorems 1 and 2.

### 3.3 Statistical Guarantee for $\hat{\beta}$ and $\hat{\Theta}$

We utilize the profiling procedure to show that  $\hat{\beta}$  and  $\hat{\Theta}$  converge to the true values. For  $\hat{\beta}$ , we show that  $\hat{\beta}$  approaches  $\beta_0$  when  $\Theta$  in the loss function is fixed at  $\hat{\Theta}$ . To reach this result, we first derive an upper bound of  $\|\partial \mathcal{L}(\hat{\Theta}, \beta_0) / \partial \beta\|_{\infty}$  in terms of  $d_{E\mathbf{X}}$  and  $d_{\mathbf{W}}$ . The upper bound vanishes as long as  $d_{E\mathbf{X}}$  and  $d_{\mathbf{W}}$  grow slower than  $mn$ . Here because we consider  $\hat{\Theta}$  instead of  $\Theta_0$  in the loss function, the upper bound is slightly larger than the standard result in high dimensional generalized linear models, where only a sparse high dimensional parameter is of interest. Furthermore, to assess the convexity of the objective function, we start with verifying that a type of the restricted eigenvalue condition (e.g. Bickel et al. (2009); Van De Geer et al. (2009)) is satisfied on a high sparsity (low  $L_1/L_2$ ) and low spikiness (low  $L_{\infty}/L_2$ ) parameter set, which contains our true parameter that

satisfies  $\|\beta\|_0 = s$  and  $\|\beta\|_\infty \leq a$ . To do that, we carefully define a set  $\mathbb{B}(D)$  (Section C), which has the following two properties: (1) the covering number of  $\mathbb{B}(D)$ , that is the number of spherical balls of a given size that cover  $\mathbb{B}(D)$ , grows slower than  $O\{\exp(mn)\}$ ; (2) its union over all values of  $D$  covers the entire space of  $\beta$ . In such  $\mathbb{B}(D)$ , we show that the probability of  $|\beta^T \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\beta - \beta^T E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\}\beta| \rightarrow 0$  at given  $\tilde{\Theta}, \tilde{\beta}$  for any  $\beta$  approaches one when  $mn \rightarrow \infty$ . In addition, we show that  $\beta^T E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\}\beta > 0$ , which leads to the result that  $\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})$  is asymptotically strictly positive definite as described in Lemma A.10. This establishes that  $\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})$  indeed satisfies the restricted eigenvalue condition. This, together with the boundedness of  $\|\partial\mathcal{L}(\hat{\Theta}, \beta_0)/\partial\beta\|_\infty$ , leads to the convergence of  $\|\hat{\beta} - \beta_0\|_2$ . We use a similar procedure to derive the upper bound for  $\|\hat{\Theta} - \Theta_0\|_F$  when  $\beta$  in the loss function is fixed at  $\hat{\beta}$ . Specifically, to study the convexity of the objective function, we verify that the restricted eigenvalue property is satisfied on a low rank (low nuclear norm/Frobenius norm ratio) and low spikiness (low  $L_{\max}/L_2$ ) set, which contains the true parameter that satisfies  $\text{rank}(\Theta) = r$  and  $\|\Theta\|_{\max} \leq a$ .

**Lemma 1** *Assume Conditions (C1), (C4) and (C5) to hold. Then there is a constant  $\omega > 0$  so that*

$$\left\| \frac{\partial\mathcal{L}(\hat{\Theta}, \beta_0)}{\partial\beta} \right\|_\infty \leq \sqrt{\frac{\omega \log\{\max(p, mn)\}}{mn}} + \frac{d_{\text{EX}}}{mn} + \frac{2ad\mathbf{w}}{mn}$$

with probability at least  $1 - 2\{\max(p, mn)\}^{-1}$ .

Lemma 1 is a direct consequence of Lemmas A.5 and A.6 in the Appendix. Here the term  $\sqrt{\omega \log\{\max(p, mn)\}}/(mn) + (mn)^{-1}d_{\text{EX}}$  is the order of  $\|\partial\mathcal{L}(\Theta_0, \beta_0)/\partial\beta\|_\infty$ , while the remaining term on the right hand side represents the order of the error  $\|\partial\mathcal{L}(\Theta_0, \beta_0)/\partial\beta - \partial\mathcal{L}(\hat{\Theta}, \beta_0)/\partial\beta\|_\infty$ .

**Theorem 1** *Assume Conditions (C1)–(C6) hold. Let*

$$\lambda_\beta \geq 2\sqrt{\frac{\omega \log\{\max(p, mn)\}}{mn}} + \frac{2d_{\text{EX}}}{mn} + \frac{4ad\mathbf{w}}{mn}.$$

Then

$$\begin{aligned} \|\hat{\beta} - \beta_0\|_2 &\leq \max \left( \left[ \frac{10\sigma_{1F}}{\alpha_{0\beta}\sqrt{mn}} + \left\{ \frac{3\lambda_\beta\sqrt{s}}{\alpha_{0\beta}} \right\}^2 \right]^{1/2} + 3\frac{\lambda_\beta\sqrt{s}}{\alpha_{0\beta}}, \right. \\ &\quad \left. 8a\sqrt{s}\gamma\sqrt{\frac{\log\{\max(p, mn)\}}{mn}}, (4\alpha_{0\beta})^{-1/2} \left\{ 2\sigma_{1F}^2 \frac{\log\{\max(p, mn)\}}{mn} \right\}^{1/4} \right) \end{aligned} \quad (2)$$

with probability at least  $1 - 4\max(p, mn)^{-1} - 2(mn)^{-1} - 2\{\max(p, mn)\}^{-C}$  for some positive constant  $C$ , where  $\sigma_{1F} = 32c_0^2a^2(d_{H_2} + d_{S_2}^2)$ , and

$$\alpha_{0\beta} \equiv \alpha_{\min}(E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\hat{\Theta}, \beta^*)\mathbf{X}_{ij}\mathbf{X}_{ij}^T|Y_{ij}\} - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\hat{\Theta}, \beta^*)\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2}])/4,$$

where  $\alpha_{\min}(\mathbf{M})$  here and throughout the text is the minimal eigenvalue of a given matrix  $\mathbf{M}$ .



The proof of Theorem 1 is lengthy and is detailed in Appendix C and D. The convergence of  $\widehat{\boldsymbol{\beta}}$  depends on many things, including the growing rate of  $p, m, n$ , the sparseness parameter  $s$ , the bounds  $d_{H_2}, d_{S_2}, d_{E\mathbf{X}}, d_{\mathbf{W}}, c_0, a, \lambda_{\boldsymbol{\beta}}$  and the missingness related quantity  $\alpha_{0\boldsymbol{\beta}}$ . If  $s$  and  $\alpha_{0\boldsymbol{\beta}}$  are finite,  $d_{\mathbf{W}}$  and  $d_{E\mathbf{X}}$  grow slower than  $mn$ , and  $\log(p) = o(mn)$ , together with the condition  $c_0^2 a^2 (d_{H_2} + d_{S_2}^2) \sqrt{\log\{\max(p, mn)\}} / (mn) \rightarrow 0$  in Condition (C6), then the upper bound of  $\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0\|_2$  in Theorem 1 will go to zero as  $mn$  increases. It is worth mentioning that  $\alpha_{0\boldsymbol{\beta}}$  always appears in the denominator of the upper bound of  $\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0\|_2$ . It indicates that although  $\text{pr}(R_{ij} = 1|Y_{ij})$  cannot be zero, it is allowed to go to zero at a certain rate. In fact, as long as its vanishing speed is sufficiently slow so that  $\lambda_{\boldsymbol{\beta}}\sqrt{s}/\alpha_{0\boldsymbol{\beta}}$  and  $\sigma_{1F}[\log\{\max(p, mn)\}/(mn)]^{1/2}/\alpha_{0\boldsymbol{\beta}}$  converge to zero, then the upper bound of  $\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0\|_2$  will still go to zero.

It is worth mentioning that, for the identifiability of the parameters, we do not allow a dense covariate matrix in the model. To see that, our Condition (C5) requires  $\|\mathbf{X}_{ij}\|_1 \leq c_0$ , where  $c_0$  is a quantity that determines the growing rate of  $\sigma_{1F} \equiv 32c_0^2 a^2 (d_{H_2} + d_{S_2}^2)$ . As shown in Theorem 1, the upper bound of the estimation error satisfies (2) in probability. Thus, for  $\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0\|_2$  to vanish to zero, we would want  $\sigma_{1F}$  to grow slower than  $\sqrt{(mn)/\log\{\max(p, mn)\}}$ . Suppose we have bounded  $d_{H_2}$  and  $d_{S_2}$ , which is the case for logistic regression, then  $c_0^2$  cannot grow faster than  $\sqrt{(mn)/\log\{\max(p, mn)\}}$ . This prevents us from having a dense covariate matrix. In fact, if the covariate vector is dense, it is likely to violate the identifiability Condition (C2) that  $\sum_{k=1}^p \beta_{0k} \mathbf{X}^k$  is a sparse matrix.

In practice, it is very rare that people can collect high dimensional dense Gaussian matrices in the social media data. In the social media data, the majority of the features are categorical, such as gender, race, address, etc, which are named dictionary in the literature (See, for example, Robin et al. (2018)). Such features are then transformed to dummy variables and result in only very few nonzero values in the covariate matrices. Some features, such as address, can have a large number of categories, and lead to a series of high dimensional but highly sparse matrices  $\mathbf{X}^k$ 's formed by dummy variables. The matrix  $\mathbf{X}^k$  thus has value one only at the entry corresponding to a specific user and a specific restaurant. In all the above situations, we are likely to have a sparse matrix  $\sum_{k=1}^p \beta_{0k} \mathbf{X}^k$ .

Our method can be applied to the majority of the social media data with the dictionary as the covariates, which is a very challenging problem even under missing-at-random settings (Robin et al., 2018). The proposed technique can be very useful in analyzing social media data.

**Lemma 2** *Assume Conditions (C1)–(C6) hold. Let*

$$\lambda_{\boldsymbol{\beta}} \geq 2\sqrt{\frac{\omega \log\{\max(p, mn)\}}{mn}} + \frac{2d_{E\mathbf{X}}}{mn} + \frac{4ad_{\mathbf{W}}}{mn}.$$

*Then for some  $c_d, \gamma > 0$ ,*

$$\begin{aligned} \left\| \frac{\partial \mathcal{L}(\boldsymbol{\Theta}_0, \widehat{\boldsymbol{\beta}})}{\partial \boldsymbol{\Theta}} \right\|_{op} &\leq c_d \sqrt{\frac{d \log(d)}{mn}} + (2d_{H_2} + 2d_{S_2}^2) \max \left( \left[ \frac{10\sigma_{1F}}{\alpha_{0\boldsymbol{\beta}}\sqrt{mn}} + \left\{ \frac{3\lambda_{\boldsymbol{\beta}}\sqrt{s}}{\alpha_{0\boldsymbol{\beta}}} \right\}^2 \right]^{1/2} \right. \\ &\quad \left. + \frac{3\lambda_{\boldsymbol{\beta}}\sqrt{s}}{\alpha_{0\boldsymbol{\beta}}}, 8a\sqrt{s}\gamma \sqrt{\frac{\log\{\max(p, mn)\}}{mn}} \right), \end{aligned}$$

$$(4\alpha_{0\beta})^{-1/2} \left\{ \frac{2\sigma_{1F}^2 \log\{\max(p, mn)\}}{mn} \right\}^{1/4}$$

with probability at least  $1 - 1/d$ .

This lemma is a direct consequence of Lemmas A.13 and A.14 in the Appendix. Here  $c_d \sqrt{d \log(d)/(mn)}$  on the right hand side captures the order of  $\|\partial \mathcal{L}(\Theta_0, \beta_0)/\partial \Theta\|_\infty$  from Lemma A.13, and the remaining terms on the right hand side describe the order of  $\|\partial \mathcal{L}(\Theta_0, \beta_0)/\partial \Theta - \partial \mathcal{L}(\Theta_0, \hat{\beta})/\partial \Theta\|_\infty$ . We further show that the estimator  $\hat{\Theta}$  described in (1) converges to the true parameter value in probability as well.

**Theorem 2** *Assume Conditions (C1)–(C6) hold. Let*

$$\begin{aligned} \lambda_{\Theta} \geq & 2c_d \sqrt{\frac{d \log(d)}{mn}} + 2(2d_{H_2} + 2d_{S_2}^2) \max \left( \left[ \frac{10\sigma_{1F}}{\alpha_{0\beta} \sqrt{mn}} + \left\{ \frac{3\lambda_{\beta} \sqrt{s}}{\alpha_{0\beta}} \right\}^2 \right]^{1/2} \right. \\ & \left. + \frac{3\lambda_{\beta} \sqrt{s}}{\alpha_{0\beta}}, 8a\sqrt{s}\gamma \sqrt{\frac{\log\{\max(p, mn)\}}{mn}}, (4\alpha_{0\beta})^{-1/2} \left\{ \frac{2\sigma_{1F}^2 \log\{\max(p, mn)\}}{mn} \right\}^{1/4} \right). \end{aligned}$$

Then

$$\begin{aligned} \|\hat{\Theta} - \Theta_0\|_F \leq & \max \left( \left[ \frac{137\sigma_{dF}}{32\alpha_{0\Theta} \sqrt{mn}} + \frac{36\lambda_{\Theta}^2 r}{\alpha_{0\Theta}^2} \right]^{1/2} + \frac{6\lambda_{\Theta} \sqrt{r}}{\alpha_{0\Theta}}, \right. \\ & \left. 16a\sqrt{r}\nu \sqrt{\frac{d \log(d)}{mn}}, (4\alpha_{0\Theta})^{-1/2} \left\{ \frac{2\sigma_{dF}^2 d \log(d)}{mn} \right\}^{1/4} \right) \end{aligned}$$

with probability at least  $1 - \exp\{-C d \log(d)\} - 2 \exp\{-d \log(d)\} - d^{-1}$ , where  $C$  is a constant. Here  $\sigma_{dF} \equiv 32a^2(d_{H_2} + d_{S_2}^2)$  and

$$\begin{aligned} \alpha_{0\Theta} \equiv & \alpha_{\min} \left( (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \Theta^*, \hat{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^\top | Y_{ij}\} \right. \\ & \left. - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \Theta^*, \hat{\beta}) \text{vec}(\mathbf{z}_{ij}) | Y_{ij}\}^{\otimes 2}] \right) / 4. \end{aligned}$$

We provide the proof of Theorem 2 in Appendix E and F. The convergence of  $\hat{\Theta}$  depends on the order of  $m, n, r, p$ . Assume that  $\Theta_0$  is a matrix with finite rank,  $d \log(d)/\log\{\max(p, mn)\} \rightarrow \infty$ , and  $d \log(d)/(mn) \rightarrow 0$ . Combining with the order of  $\|\partial \mathcal{L}(\Theta_0, \hat{\beta})/\partial \Theta\|_{op}$  in Lemma 2, the upper bound of  $\|\hat{\Theta} - \Theta_0\|_F$  is of order  $O_p\{d \log(d)/(mn)^{1/4}\}$ . Similar to Theorem 1, we can see that  $\alpha_{0\Theta}$  always appears in the denominator in the upper bound of  $\|\hat{\Theta} - \Theta_0\|_2$ . This implies  $\text{pr}(R_{ij} = 1 | Y_{ij})$  is allowed to go to zero as long as the vanishing speed is sufficiently slow so that  $\lambda_{\Theta} \sqrt{r}/\alpha_{0\Theta}$  and  $\sigma_{dF}\{d \log(d)/(mn)\}^{1/2}/\alpha_{0\Theta}$  converge to zero.

Due to the high dimensionality of  $\beta$  and  $\Theta$ , the convergence of  $\hat{\beta}$  and  $\hat{\Theta}$  is difficult to achieve. The sample version of the loss function can be non-convex, which may lead to non-unique solution of the estimation procedure. To show statistical convergence, we show

that the second derivative of the loss function with respect to the parameters satisfies the restricted eigenvalue conditions as shown in Lemmas A.10 and A.19 for  $\beta$  and  $\Theta$ , respectively. These restricted eigenvalue conditions guarantee the convergence of our estimators in the feasible set. Unlike standard results, the convergence of  $\widehat{\Theta}$  depends on that of  $\widehat{\beta}$  and vice versa, and hence we require  $\log\{\max(p, mn)\} = o_p(mn)$  and  $d\log(d) = o_p(mn)$  to achieve the consistency of both  $\widehat{\beta}$  and  $\widehat{\Theta}$ . The detailed theoretical derivations are in the Appendix. We first establish a series of necessary lemmas in Sections C and E of the Appendix, and then provide the detailed proofs in Sections D and F of the Appendix.

#### 4. Computational Algorithm and Convergence

We now describe the computational algorithm and show its convergence rate theoretically. In the computational aspect, we treat the two parameters  $\Theta_0$  and  $\beta_0$  differently. In updating  $\beta$  we use the proximal gradient algorithm, and in updating  $\Theta$  we use soft-impute. We name the combination of the proximal gradient algorithm and the soft-impute treatment as the Soft Impute Proximal Gradient (SIPG) algorithm. In computing the integrals, we can view the integrations as mean and approximate it using sample averages across different  $\mathbf{X}_{ij}$ 's when this brings computational gain.

**Remark 1** Note that  $\partial\mathcal{L}(\Theta, \beta)/\partial\beta - \partial\mathcal{L}(\Theta, \beta')/\partial\beta = \langle \partial\mathcal{L}(\Theta, \beta'')^2/\partial\beta\partial\beta^T, \beta - \beta' \rangle$ , by Conditions (C5) and (C6), where  $\beta''$  is a point on the line connecting  $\beta$  and  $\beta'$ . Therefore,  $\|\partial\mathcal{L}(\Theta, \beta)/\partial\beta - \partial\mathcal{L}(\Theta, \beta')/\partial\beta\|_2 \leq \sigma_\beta\|\beta - \beta'\|_2$  with  $\sigma_\beta \equiv \{2c_0^2(d_{H_2} + d_{S_2}^2)\}$ . We say that  $\mathcal{L}(\Theta, \beta)$  is  $\sigma_\beta$  smooth with respect to  $\beta$ . Similarity let  $\Theta$  and  $\Theta'$  be two  $m \times n$  matrices, and  $\Theta'_{lk} = \Theta_{lk}$  if  $(l, k) \neq (i, j)$ , and  $\Theta'_{ij} \neq \Theta_{ij}$ . It is easy to show that

$$|\partial\mathcal{L}(\Theta, \beta)/\partial\Theta_{ij} - \partial\mathcal{L}(\Theta', \beta)/\partial\Theta_{ij}| \leq \sigma_\Theta|\Theta_{ij} - \Theta'_{ij}|,$$

where  $\sigma_\Theta \equiv 2(d_{H_2} + d_{S_2}^2)$ . We say that  $\mathcal{L}(\Theta, \beta)$  is  $\sigma_\Theta$  smooth with respect to  $\Theta$ . Similarly we can show that

$$\|\partial\mathcal{L}(\Theta, \beta)/\partial\beta - \partial\mathcal{L}(\Theta', \beta)/\partial\beta\|_2 \leq \sigma_{\beta\Theta}\|\Theta - \Theta'\|_F,$$

for  $\sigma_{\beta\Theta} \equiv 2c_0(d_{H_2} + d_{S_2}^2)$ . We say that  $\partial\mathcal{L}(\Theta, \beta)/\partial\beta$  is  $\sigma_{\beta\Theta}$  Lipschitz continuous.

Using  $\beta^t, \Theta^t$  to denote the corresponding estimators at the  $t$ th iteration, and assuming in the  $t$ th iteration  $\beta^{t-1}, \Theta^{t-1}$  are given, we estimate  $\beta^t$  using the proximal gradient descent. That is, we obtain

$$\beta^t = \operatorname{argmin}_\beta \frac{1}{2}\|\beta - \beta^{t-1} + \eta \frac{\partial\mathcal{L}(\Theta^{t-1}, \beta^{t-1})}{\partial\beta}\|_2^2 + \eta\lambda_\beta\|\beta\|_1$$

through letting

$$\beta^t = \rho_{\eta\lambda_\beta} \left( \beta^{t-1} - \eta \frac{\partial\mathcal{L}\{\Theta^{t-1}, \beta^{t-1}\}}{\partial\beta} \right),$$

where  $\rho_a$  is the component-wise soft thresholding operator so that the  $i$ th element of  $\rho_a(\mathbf{x})$  is  $x_i - a$  if  $x_i > a$ ,  $x_i + a$  if  $x_i < -a$  or 0 if  $|x_i| < a$ . Here  $\eta$  is a step size. To update  $\Theta$ , we

consider

$$\Theta^t = \operatorname{argmin}_{\Theta} \frac{1}{2} \left\| \Theta - \Theta^{t-1} + \eta_1 \frac{\partial \mathcal{L}(\Theta^{t-1}, \beta^t)}{\partial \Theta} \right\|_F^2 + \eta_1 \lambda_{\Theta} \|\Theta\|_*,$$

where  $\eta_1$  is the step size, and use Soft-Impute (Mazumder et al., 2010) algorithm to implement the optimization. Specifically, let  $\mathbf{UDV}^T$  be the singular value decomposition of  $\Theta^{t-1} - \eta_1 \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \Theta$ , we obtain

$$\Theta^t = \mathbf{U} \operatorname{Soft}_{\eta_1 \lambda_{\Theta}}(\mathbf{D}) \mathbf{V}^T,$$

where  $\operatorname{Soft}_{\eta_1 \lambda_{\Theta}}(\cdot)$  operates element-wise on the diagonal matrix  $\mathbf{D}$  by replacing  $D_{ii}$  with  $(D_{ii} - \eta_1 \lambda_{\Theta})_+$ .

Let  $F(\Theta, \beta) = \mathcal{L}(\Theta, \beta) - L + \lambda_{\Theta} \|\Theta\|_* + \lambda_{\beta} \|\beta\|_1$  for  $|L| < \infty$  such that  $\mathcal{L}(\Theta, \beta) > L$ . In Theorem 3, we establish the convergence of the proposed SIPG algorithm.

To prepare for Theorem 3, define

$$\begin{aligned} g_{\beta}(\Theta, \beta, Q) &= \max_{\|\tilde{\beta}\| \leq Q} \langle \partial \mathcal{L}(\Theta, \beta) / \partial \beta, \beta - \tilde{\beta} \rangle + \lambda_{\beta} (\|\beta\|_1 - \|\tilde{\beta}\|_1), \\ g_{\Theta}(\Theta, \beta, R) &= \max_{\|\tilde{\Theta}\|_* \leq R} \langle \partial \mathcal{L}(\Theta, \beta) / \partial \Theta, \Theta - \tilde{\Theta} \rangle + \lambda_{\Theta} (\|\Theta\|_* - \|\tilde{\Theta}\|_*). \end{aligned}$$

Let  $Q^t = F(\Theta^t, \beta^t) / \lambda_{\beta}$ ,  $R^t = F(\Theta^t, \beta^{t+1}) / \lambda_{\Theta}$ ,  $\sigma_{\beta} = C_1(2d_{H_2} + 2d_{S_2}^2)c_0^2$  for  $C_1 > 1$ ,  $\sigma_{\Theta} = 2d_{H_2} + 2d_{S_2}^2$ .

**Theorem 3** Assume that for all  $\beta, \tilde{\beta}$  that satisfy  $\|\beta\|_{\infty} \leq 2a$ ,  $\|\tilde{\beta}\|_{\infty} \leq a$  and all  $\Theta, \tilde{\Theta}$  that satisfy  $\|\Theta\|_{\max} \leq 2a$ ,  $\|\tilde{\Theta}\|_{\max} \leq a$ ,

$$\begin{aligned} \beta^T E\{\mathbf{F}_{\beta}(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta})\} \beta &\geq 4c_0 \sqrt{(d_{H_2} + d_{S_2}^2) \frac{\log\{\max(p, mn)\}}{mn}} \|\beta\|_2^2, \\ \operatorname{vec}(\Theta)^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta})\} \operatorname{vec}(\Theta) &\geq 4 \sqrt{(d_{H_2} + d_{S_2}^2) \frac{\log(mn)}{mn}} \|\Theta\|_F^2, \\ E \left( \left\{ \begin{array}{c} \beta \\ \operatorname{vec}(\Theta) \end{array} \right\}^T \frac{\partial^2 \mathcal{L}(\Theta, \beta)}{\partial [\{\beta^T, \operatorname{vec}(\Theta)^T\}^T]^{\otimes 2}} \left\{ \begin{array}{c} \beta \\ \operatorname{vec}(\Theta) \end{array} \right\} \right) &\geq 4(c_0 + 1) \sqrt{(d_{H_2} + d_{S_2}^2) \frac{\log(mn)}{mn}} \left\| \left\{ \begin{array}{c} \beta \\ \operatorname{vec}(\Theta) \end{array} \right\} \right\|_2^2 \end{aligned} \quad (3)$$

where the expectations are taken over  $\mathbf{X}, \mathbf{Y}$ . Select  $0 < \eta < 1/\sigma_{\beta}$  and select  $0 < \eta_1 < 1/\sigma_{\Theta}$ . Then when

$$\sum_{t=0}^{T-1} 1/C(t) \geq \left\{ \frac{1}{\epsilon} - \frac{1}{F(\Theta^0, \beta^0) - F(\hat{\Theta}, \hat{\beta})} \right\} \quad (4)$$

for  $C(t) = \max\{32(R^{t-1})^2/\eta_1 + 16\sigma_{\beta\Theta}^2(Q^t)^2/(1/\eta_1 - \sigma_{\Theta}), 32(Q^t)^2/\eta\}$ , we have

$$F(\Theta^T, \beta^T) - F(\hat{\Theta}, \hat{\beta}) \leq \epsilon$$

with probability at least  $1 - 6(mn + p)^{-1}$ . Hence  $F(\Theta^T, \beta^T)$  is the  $\epsilon$ -optimal solution for (1) with probability at least  $1 - 6(mn + p)^{-1}$ .

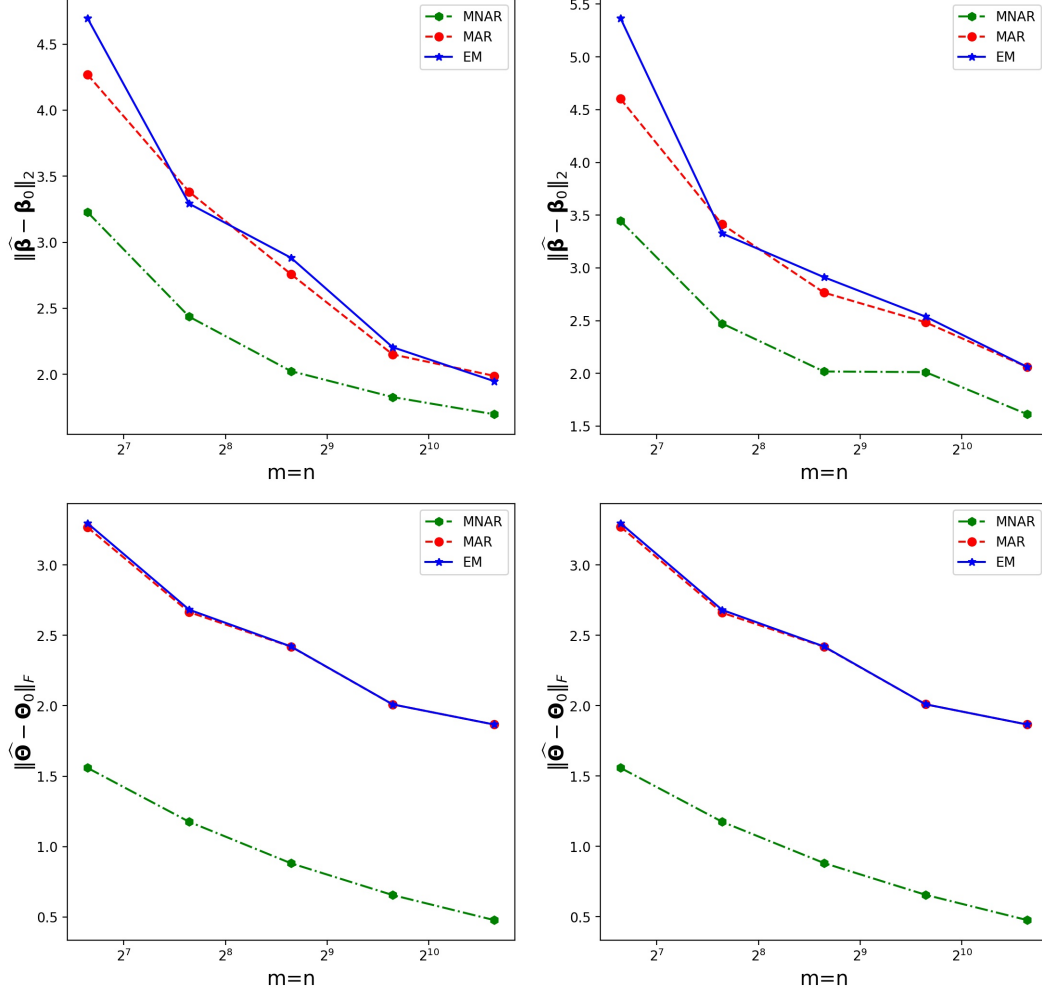


Figure 1: Linear Model : The averages of  $\|\hat{\beta} - \beta_0\|_2$  and  $\|\hat{\Theta} - \Theta_0\|_F$  of MNAR, EM and MAR over 50 simulations when  $p = 50$  (left) and  $p = 100$  (right).

The proof of Theorem 3 is given in Appendix G and H. Theorem 3 shows that when  $mn + p \rightarrow \infty$ , with probability approaching 1, the SIPG algorithm indeed converges to the optimizer described in (1) as long as sufficiently many iterations are carried out. Here the left hand side of (4) depends on the number of iterations  $T$ .

**Remark 2** In Theorem 3, (3) leads to the convexity of the objective function when the parameter search is implemented in the feasible set of  $\|\beta\|_\infty \leq a$  and  $\|\Theta\|_{\max} \leq a$  and satisfies (3). By the definition of  $Q^t, R^t$ , we have  $Q^t \leq Q^0$  and  $R^t \leq R^0$  for  $t \geq 1$ . Therefore,

$$C(t) \leq \bar{C} \equiv \max\{32(R^0)^2/\eta_1 + 16\sigma_{\beta\Theta}^2(Q^0)^2/(1/\eta_1 - \sigma_\Theta), 32(Q^0)^2/\eta\},$$

i.e.,  $\sum_{t=0}^T \{1/C(t)\} \geq T\bar{C}^{-1}$ . Thus, a sufficient condition for (4) is

$$T \geq \bar{C} \left\{ \frac{1}{\epsilon} - \frac{1}{F(\boldsymbol{\Theta}^0, \boldsymbol{\beta}^0) - F(\hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\beta}})} \right\}.$$

This indicates that regardless what is the initial value, as long as the convexity on the feasible set holds, we can always perform sufficiently many iterations to obtain the  $\epsilon$ -optimal solution with probability approaching to one. Of course, a closer initial value  $(\boldsymbol{\Theta}^0, \boldsymbol{\beta}^0)$  to the solution will require fewer iterations. We point out that our conditions in (3) are very mild and are weaker than the standard literature. In fact, because the population version of  $\mathcal{L}(\boldsymbol{\Theta}, \boldsymbol{\beta})$  is a convex function, the expectations in (3) are naturally nonnegative. In the existing literature (see, for example, the restrict eigenvalue condition in Loh and Wainwright (2012)), it is standard to further require the expectations to be strictly positive, i.e., the expectations in (3) are usually required to be larger than a positive constant. Here, because we allow the missing probability to go to 1, we relax the strictly positive condition by replacing the positive constant with a positive value that goes to zero.

## 5. Simulations

We perform three simulation studies to evaluate the MNAR method. The simulations are repeated 50 times. We first generate  $X_{ijk} \sim \text{Bernoulli}(q)$  for  $k = 1, \dots, p$  independently, where  $q = 0.2$ . Thus,  $g(\mathbf{x}_{ij}) = \prod_{k=1}^p q^{x_{ijk}} (1-q)^{1-x_{ijk}}$ . Then we generate  $Y_{ij}$  from a linear model with mean  $\boldsymbol{\Theta}_{0ij} + \boldsymbol{\beta}_0^T \mathbf{X}_{ij}$  and standard deviation 5; We design  $\boldsymbol{\Theta}_0$  to be a rank 5 matrix with singular values (10, 1.8, 1.6, 1.4, 1.2), and  $\boldsymbol{\beta}_0 = (1, 0, 2, 0, -3, -4, 5, 0, \dots, 0)^T$ . We set  $m = n$ , and vary  $m, n$  from 100 to 1600. Furthermore, we generate  $R_{ij}$  from  $\Pr(R_{ij} = 1 | Y_{ij}) = \text{expit}(Y_{ij} - \bar{Y} - D)$ , where  $\bar{Y} = \sum_{i,j} Y_{ij} / (mn)$  and  $D$  is chosen to achieve 90% missingness in the data. As specified in Theorems 1 and 2, we select  $\lambda_{\boldsymbol{\beta}} = C_{\boldsymbol{\beta}} \sqrt{\log\{\max(p, mn)\} / mn}$  and  $\lambda_{\boldsymbol{\Theta}} = C_{\boldsymbol{\Theta}} \max\left\{\sqrt{\log(d)/d}, \log\{\max(p, mn)\}^{1/4} \sqrt{d}\right\}$ , where  $C_{\boldsymbol{\beta}}$  and  $C_{\boldsymbol{\Theta}}$  are constants chosen to achieve similar sparseness of the estimators across all situations.

In Figure 1, we compare our method with the likelihood method which uses only the observed data and assume missing at random (referred to as MAR method), and the expectation-maximization (EM) method which imputes the missing outcomes by their expected values based on the previous estimators. The results show the MNAR method outperforms the EM and MAR methods consistently over all settings.

In the second simulation, we generate the response  $Y_{ij}$ 's from a logistic model with mean  $\text{expit}(\boldsymbol{\Theta}_{0ij} + \boldsymbol{\beta}_0^T \mathbf{X}_{ij})$ , and we generate  $R_{ij}$  from  $\Pr(R_{ij} = 1 | Y_{ij}) = 0.19I(Y_{ij} = 1) + 0.01I(Y_{ij} = 0)$ , so the missing probability is around 0.90 marginally. Figure 2 shows that the MNAR method also outperforms the EM and MAR methods consistently under the logistic model.

Furthermore, we conduct additional simulations when the covariates are row (column)-specific, meaning that the covariates values are the same for the observations in the same column (row). We adopt the same simulation setting as the ones used in Figure 1, except that we change the covariate structure, where the first half covariates are row-specific while

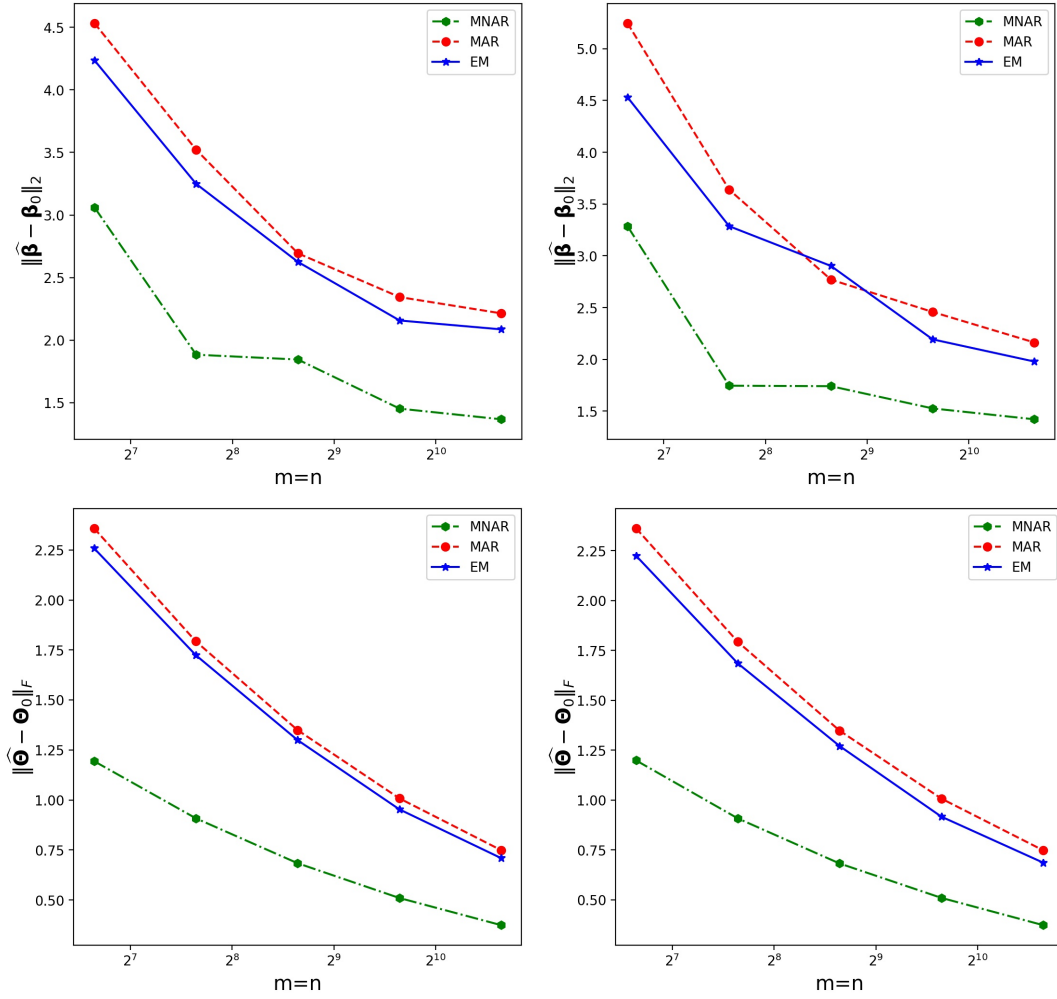


Figure 2: Logistic Model: The averages of  $\|\hat{\beta} - \beta_0\|_2$  and  $\|\hat{\Theta} - \Theta_0\|_F$  of MNAR, EM and MAR over 50 simulations when  $p = 50$  (left) and  $p = 100$  (right).

the last half covariates are column-specific. As shown in Figure 3, MNAR still outperforms the other two methods under this setting.

Lastly, to evaluate the performance of MNAR in larger sample settings, we generate data when  $p = 2$ ,  $\beta_0 = (1, -2)^T$  and vary  $m, n$  from 100 to 12800 with  $m = n$ . We compare the averages of  $\|\hat{\beta} - \beta_0\|_2$  and  $\|\hat{\Theta} - \Theta_0\|_F$  of MNAR, EM and MAR over 50 simulations in Figure 4. The results show that the MNAR method outperforms the EM and MAR methods in this setting. The experiment also demonstrates the ability of MNAR to handle large data sets.

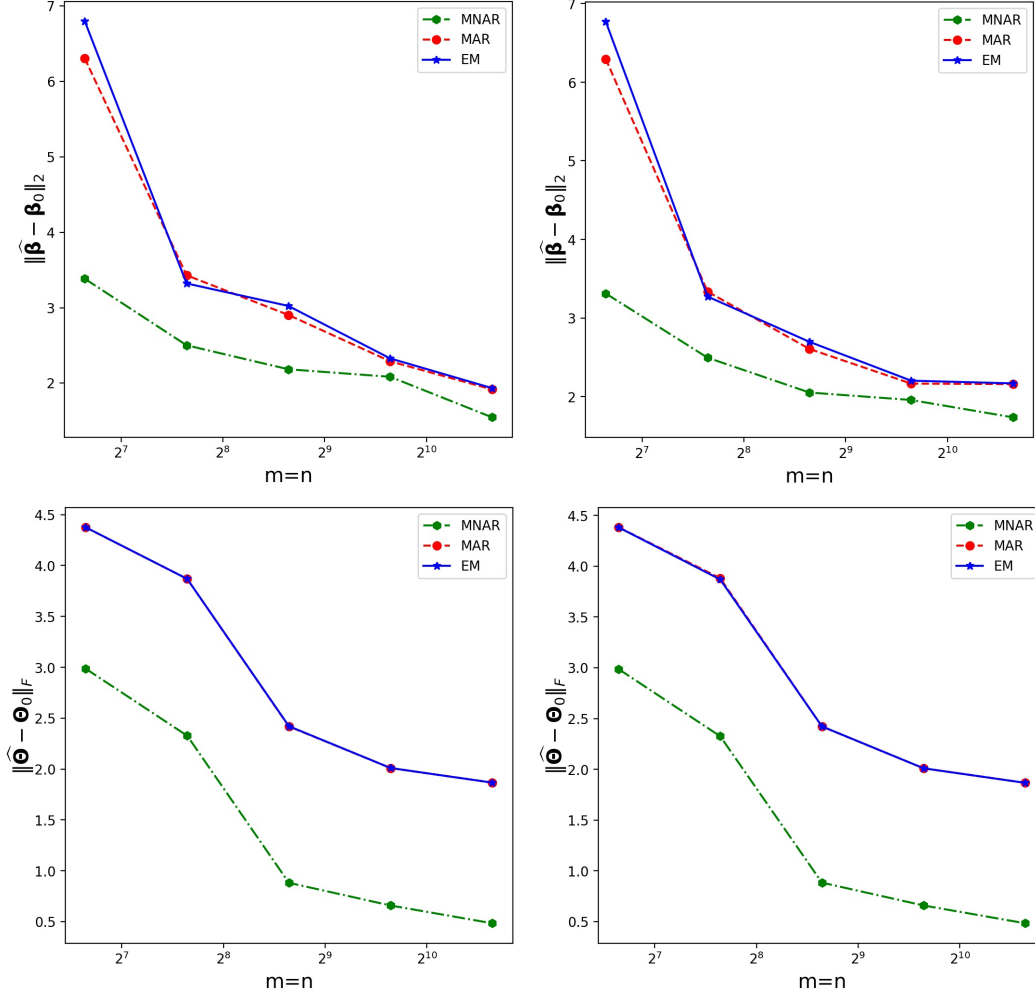


Figure 3: Row (column) specific covariates: The averages of  $\|\hat{\beta} - \beta_0\|_2$  and  $\|\hat{\Theta} - \Theta_0\|_F$  of MNAR, EM and MAR over 50 simulations when  $p = 50$  (left) and  $p = 100$  (right).

## 6. Real Data Analysis

We evaluate the performance of the MNAR, MAR, and EM methods on the real data from Yelp. Furthermore, we compare MNAR with a benchmark weighted collaborative filtering algorithm on the MovieLens (<https://grouplens.org/datasets/movielens/>) and the Yelp data.

### 6.1 Yelp Data Analysis

We apply the proposed method to analyze the data from Yelp, and the data set is available at <https://www.yelp.com/dataset/documentation/main>. The full data set is huge, which



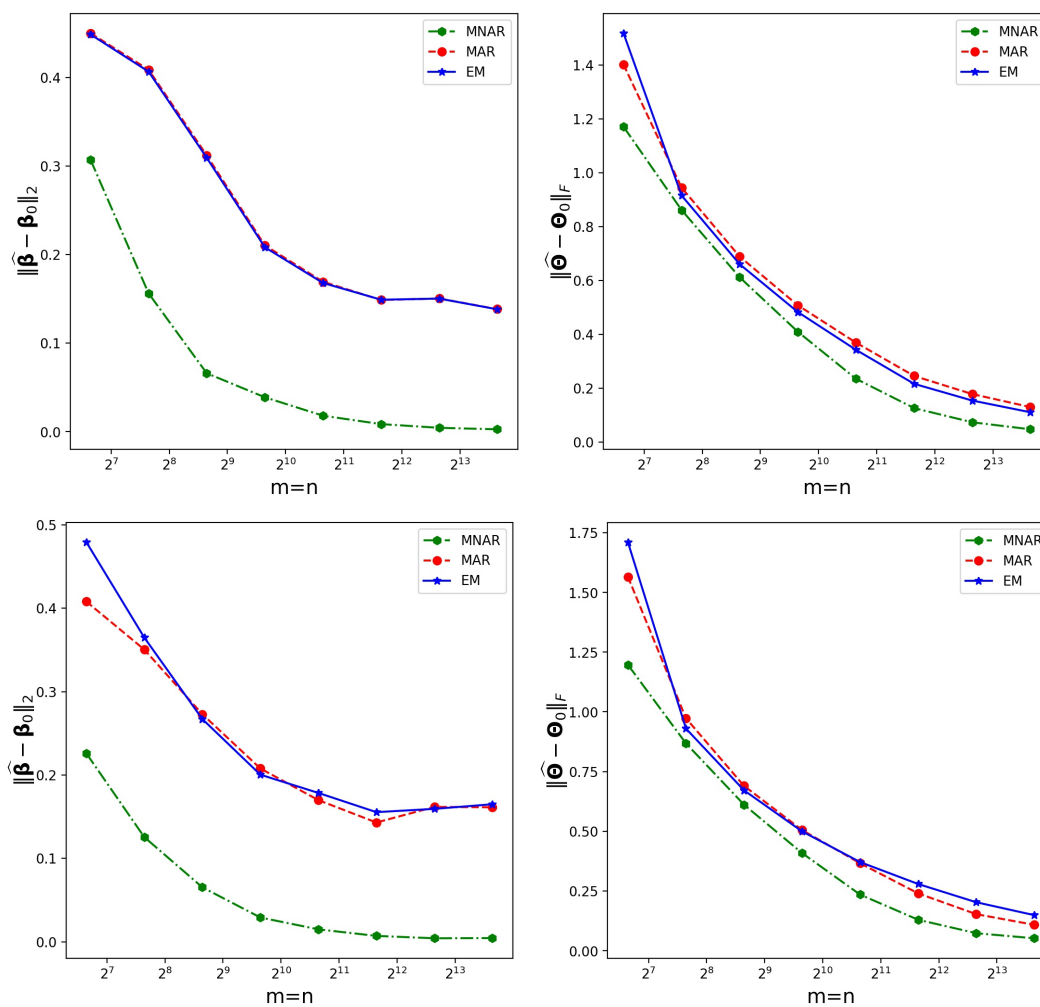


Figure 4: Large matrix size: The averages of  $\|\hat{\beta} - \beta_0\|_2$  and  $\|\hat{\Theta} - \Theta_0\|_F$  of MNAR, EM and MAR over 50 simulations under linear (upper) and logistic models (lower) when  $p = 2$ .

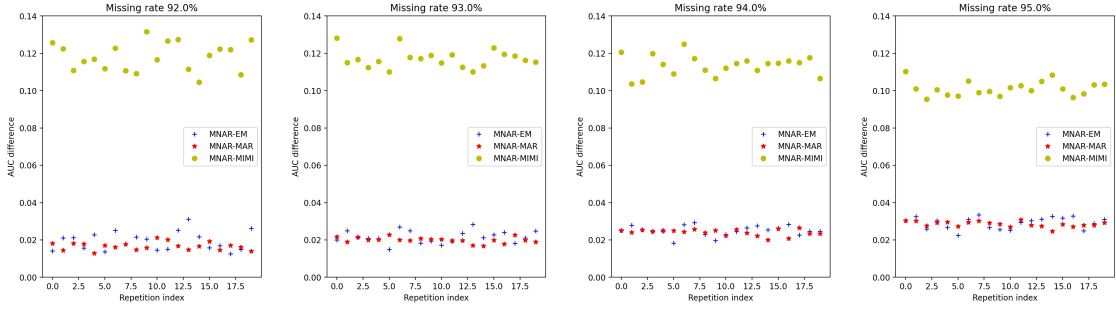


Figure 5: The AUC differences between MNAR and the other three models for 20 repetitions under different missing rates.

includes 8635403 reviews on 160585 businesses from 8 metropolitan areas. In combination with a large number of covariates, it leads to prohibitive computation. Thus, we first select 500 restaurants which received most ratings in Las Vegas. Then we choose 1000 users who gave most ratings for these 500 restaurants. The resulting  $\mathbf{Y}$  matrix contains the reviews from 1000 ( $n = 1000$ ) customers at 500 ( $m = 500$ ) restaurants with 90.9% missingness rate. In addition, each  $Y_{ij}$ 's corresponds to a 22 ( $p = 22$ ) dimensional covariates vector  $\mathbf{X}_{ij}$ , which includes baseline features of the  $i$ th customer and the  $j$ th restaurant such as the restaurant star, the open date, the customer review count, etc. We further standardize these covariates by subtracting the mean and dividing the standard deviation. We dichotomize the responses so that  $Y_{ij} = 1$  if the review score from the  $i$ th customer at the  $j$ th restaurant is greater than 3.5. Furthermore, we introduce an evaluation procedure specifically for this missing not at random data set as follows. We first remove the observed  $Y_{ij}$ 's with probabilities  $p_1$  and  $p_0$  for  $Y_{ij} = 1$  and  $Y_{ij} = 0$ , respectively, where  $p_0$  and  $p_1$  with  $p_0 \neq p_1$  are chosen so that an additional  $\alpha 100\%$  missingness is introduced into the data. Then we perform the proposed method on the remaining data, and use the estimation results to predict the removed entries. We perform the procedure 20 times, and compare the area under the receiver operating characteristic curve (AUC) of the MNAR, MAR, EM and MIMI methods proposed by Robin et al. (2020) under the settings with  $\alpha = 0.011, 0.021, 0.031, 0.041$ . In the implementation, we use the Monte Carlo method to approximate the integration in the loss function, while the distribution of  $\mathbf{X}_{ij}$  is estimated empirically. Because  $\beta^T \mathbf{X}_{ij}$  contributes to the conditional distribution of  $Y_{ij}$ , we re-estimate the distribution of  $\beta^T \mathbf{X}_{ij}$  when a new  $\beta$  is obtained. At the  $t$ th iteration, we sample  $M$  copies of  $\beta^{tT} \mathbf{X}_{ij}$  from the empirical distribution, and approximate the integration as

$$(mn)^{-1} \sum_{s=1}^m \sum_{u=1}^n f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \beta^{tT} \mathbf{X}_{su}) \approx \int f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \beta^{tT} \mathbf{X}) g(\mathbf{X}) d\mathbf{X}.$$

The penalty parameters in the four methods are selected by grid search, i.e., under each  $\alpha$ , we generate a missing matrix  $\mathbf{R}$  to tune the penalty parameters which yield largest AUC for the four methods.

We plot the 20 AUC differences between MNAR and the other three methods in Figure 5. In addition, we plot the average AUCs of the four methods over the 20 repetitions in

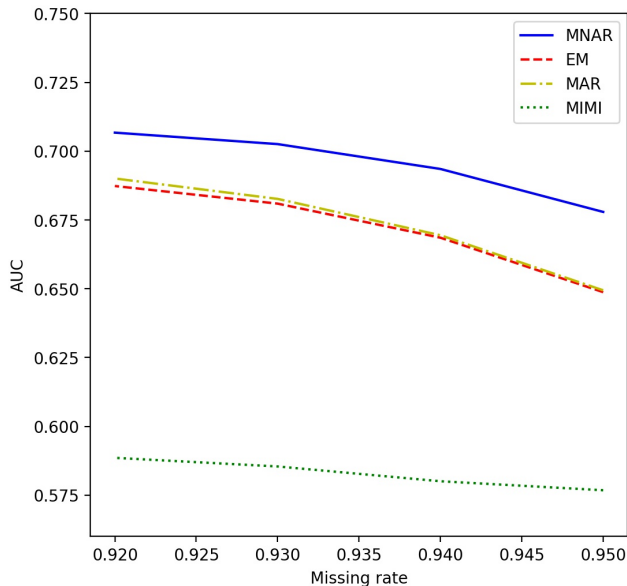


Figure 6: The means of AUC for the MNAR, EM, MAR and MIMI based on 20 repetitions under different missing rates.

Figure 6. The results show that MNAR clearly outperforms all competing methods with the largest AUC consistently across all settings.

## 6.2 Comparison with the weighted collaborative filtering method

Weighted collaborative filtering (WCF) is arguably the first approach that tackles the problem of nonignorable missing by assuming the missing data are the places where most negative feedbacks are expected to be found (Hu et al., 2008). We compare the performance of our method with that of the WCF method proposed in Hu et al. (2008) on MovieLens 1M data set, which includes one million ratings from 6040 users and 3952 movies. The data set also contains information on the age of a user and genre of a movie, which we use as the covariates in the MNAR method. There are 7 age groups and 302 different movie genres, so the covariate size is  $p = 7 \times 302 = 2114$ .

The missingness rate is over 95.7% in the full data. To evaluate the two methods, we first remove the observed  $Y_{ij}$ 's with probabilities  $p_1$  and  $p_0$  for  $Y_{ij} = 5$  and  $Y_{ij} \neq 5$  respectively, where  $p_0$  and  $p_1$  ( $p_0 \neq p_1$ ) are chosen so that an additional 1% missingness is introduced into the training data. We evaluate the performance of the two methods using the expected percentile ranking proposed by Hu et al. (2008). Specifically, after obtaining the predicted ratings of all the movies by all the users, we obtain the percentile-ranking  $r_{ui}$  of movie  $i$  in all movies (including both rated and unrated movies) for user  $u$ .

Now let  $Y_{ui}^t$  be the observed rating of movie  $i$  by user  $u$  in the testing data set. Then the expected percentile ranking is calculated as

$$\bar{r} = \frac{\sum_{ui} Y_{ui}^t r_{ui}}{\sum_{ui} Y_{ui}^t}.$$

Note that a smaller  $\bar{r}$  indicates better performance. The same procedure is repeated 20 times, and we provide the average of the 20  $\bar{r}$ 's as the final performance measure.

For the WCF method, we obtain the weights using the BM25 function (Christopher et al., 2008) with two tuning parameters  $b$ ,  $k_1$  selected as  $b = 0.75$  and  $k_1 = 2$ . This selection gives us the best performance among those recommended in Christopher et al. (2008) ( $b = 0.75$  and  $k_1 \in [1.2, 2.0]$ ). We also plot the mean of the expected percentile ranking over different choices of factors on 20 testing data sets in the left panel of Figure 7. The results show that WCF is sensitive to the selection of the number of factors with the optimal number of factors being 16 among all selections. Moreover, we plot the boxplots of the expected percentile rankings over the 20 testing data sets from the MNAR and WCF methods in the right panel of Figure 7. The results show that MNAR outperforms WCF with significantly smaller expected percentile rankings.

We also apply MNAR and WCF on the Yelp data set discussed in Section 6.1. We again plot the mean of the expected percentile ranking over different choices of factors in the left panel of Figure 8. The results show that the WCF is sensitive to the selection of the number of factors with the optimal number of factors being 4 among all selections. Moreover, we plot the boxplots of the expected percentile rankings from MNAR and WCF methods in the right panel of Figure 8. The results show that MNAR outperforms WCF with significantly smaller expected percentile rankings.

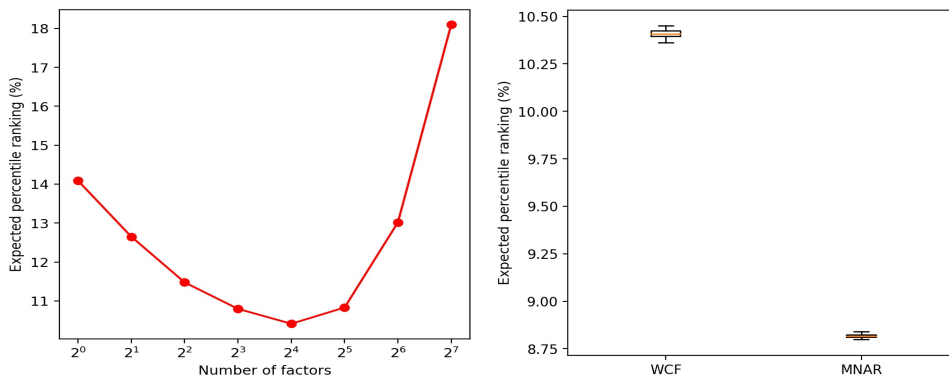


Figure 7: MovieLens data analysis: Left: The expected percentile ranking versus the number of factors from the WCF method. The expected percentile ranking decreases when the number of factors is less than 16, and starts to increase when the number of factors is greater than 16. Right: The boxplots of the expected percentile rankings from MNAR and WCF. MNAR has significantly smaller expected percentile rankings.

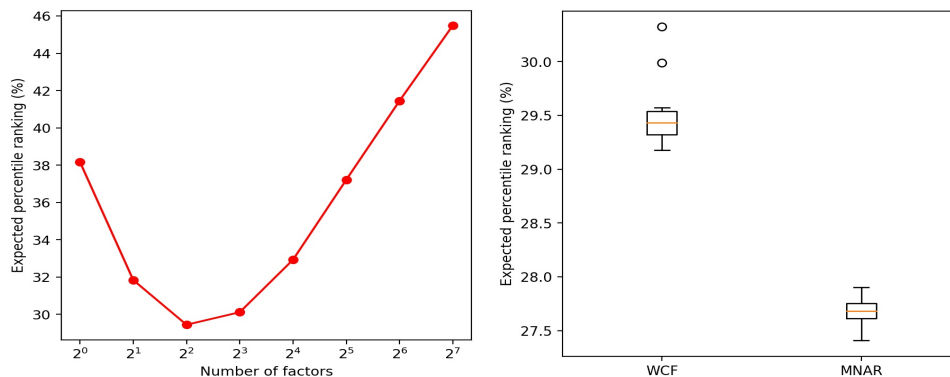


Figure 8: Yelp data analysis: Left: The expected percentile ranking versus the number of factors from the WCF method. The expected percentile ranking decreases when the number of factors is less than 4, and starts to increase when the number of factors is greater than 4. Right: The boxplots of the expected percentile rankings from MNAR and WCF methods. MNAR has significantly smaller expected percentile rankings.

## 7. Related work

Matrix completion problems have caught researchers' attention in recent years. The earlier works in this area do not consider any noise in the response (Candés and Recht, 2009; Recht, 2011), and only study the theoretical properties of perfect matrix completion. Noise issue is later taken into account in Candés and Plan (2010); Koltchinskii et al. (2011); Rohde and Tsybakov (2011); Negahban and Wainwright (2012), although they still do not include covariates or study missingness mechanism. In more recent works, nonuniform missing mechanism has also raised attention (Srebro and Salakhutdinov, 2010; Negahban and Wainwright, 2012; Klopp, 2014; Cai and Zhou, 2016; Cai et al., 2016; Bi et al., 2016; Mao et al., 2019). In addition, covariate information is taken into account. For example, Abernethy et al. (2009); Xu et al. (2013); Chiang et al. (2015); Mao et al. (2018) consider additional finite dimensional covariate effect; Zhu et al. (2016) cast the matrix completion problem in a sparse regression setting and estimate the high dimensional parameters with the conventional lasso type penalty; Robin et al. (2018) consider matrix completion with high dimensional covariate in the generalized linear model when the loss function is convex. Nevertheless, even the most sophisticated missingness feature studied in Mao et al. (2019) and the most flexible covariate structure studied in Robin et al. (2018) still assume a specific missingness mechanism model and require that the missingness does not depend on the potential response value given the covariates, i.e. they require modeling the missingness mechanism and they are limited to the framework of missing at random. These assumptions conveniently facilitate the parameter estimation, because the missingness does not affect the convexity of the loss function, and the estimation consistency can be easily achieved when the proportion of missing is not overwhelmingly large. However, this is restrictive and

is somewhat unrealistic in many customer evaluation system such as the Yelp restaurant evaluation problem considered here.

Missingness is usually classified into three categories, missing completely at random, missing at random and missing not at random. The first two classes have received extensive studies and a vast amount of literature is available. However, the third class, which is sometimes also referred to as informative or nonignorable missing, is much harder to treat and has not been well studied until very recent years. Some attempts are made, see, for example, Zhao and Shao (2015); Miao and Tchetgen Tchetgen (2016). The nonignorable missing problem has also caught attention in the collaborative filtering literature. Marlin and Zemel (2009), Hernández-Lobato et al. (2014), and Liang et al. (2016) discuss the collaborative prediction and ranking with nonignorable missing data while assuming the missing distributions are given; Hu et al. (2008) and Pan et al. (2008) tackle the problem of MNAR data in the case of implicit feedback-data; and Steck (2010) discusses how to evaluate the algorithm under MNAR settings. But none of these works directly predicts the product rating under the MNAR setting for the unknown missing distribution as we do, and none of them is in the high dimensional covariates setting as we consider here. Our work is in the third category, with the additional complexity of ultra-high dimensional covariates. To handle the nonignorable missing, we establish a penalized pseudo-likelihood framework, where despite of the goal of predicting  $Y_{ij}$  based on the covariate information  $\mathbf{X}_{ij}$ , we work with the reverse relation of  $\mathbf{X}_{ij}$  given  $Y_{ij}$ , hence bypassing the difficulties caused by the nonignorable missingness. However, in the ultra-high dimensional covariate situation, this leads to a non-convexity issue even if the original full data likelihood is convex. We address and overcome these issues caused by the ultra-high dimension by simultaneously incorporating the low-rank restriction on the baseline evaluation matrix and the sparseness assumption on the potentially high dimensional covariates. Moreover, we develop efficient computational algorithms to obtain the low-rank and sparse parameter estimators. To the best of our knowledge, our work is the first complete framework to allow the missingness in the matrix to depend on the evaluation value that is itself missing, while simultaneously considering ultra-high dimensional covariates. The penalized pseudo-likelihood strategy may also have wide application in supervised learning settings when outcomes are missing with unknown mechanisms, or more generally when the sample is biased due to various sampling issues. We aim at tackling the complex case where both the high dimensional covariate effect and matrix completion are of interest. The proposed method is easily generalizable to simpler settings when only the covariate effect is of interests or only the matrix completion with low dimensional covariate effects is considered.

## 8. Conclusion

We have considered a matrix completion problem where not only the missingness mechanism of the entries in the matrix belongs to the missing not at random context, but also the covariate information is taken into account. The diverging size of the matrix is regularized through a rank constraint enforced through penalizing its nuclear norm, and the diverging dimension of the covariates is regularized through a sparsity assumption imposed via the  $L_1$  penalization. It will be of interest to investigate if other practically justifiable constraints can be used to replace or enrich these assumptions. It will also be of interest to investigate

if better procedures can be developed if one is willing to make more concrete assumptions on the missingness mechanism. These are challenging problems but can be rewarding to study.

### **Acknowledgments**

We are grateful to the referees and the editor for valuable comments and constructive suggestions. The authors are partially supported by grants from NSF and NIH.

## Appendix

### Appendix A. The Expression of Notations in Section 3.1

To facilitate the theoretic derivation, we write out the explicit expression for some important quantities. Recall the definition of  $\mathbf{z}_{ij}$  defined in Section 3.1. Now

$$\begin{aligned} & \frac{\partial \mathcal{L}(\Theta, \beta)}{\partial \Theta} \\ = & -(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} \left[ \left\{ \frac{f_2(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})}{f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})} - \frac{\int f_2(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X}}{\int f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X}} \right\} \mathbf{z}_{ij} \right] \end{aligned}$$

and

$$\begin{aligned} & \frac{\partial^2 \mathcal{L}(\Theta, \beta)}{\partial \text{vec}(\Theta) \partial \text{vec}(\Theta)^T} \\ = & -(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} \left[ \left\{ \frac{f_{22}(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})}{f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})} - \frac{f_2^2(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})}{f^2(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})} \right. \right. \\ & \left. \left. - \frac{\int f_{22}(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X}}{\int f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X}} + \frac{\left\{ \int f_2(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X} \right\}^2}{\left\{ \int f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X} \right\}^2} \right\} \text{vec}(\mathbf{z}_{ij})^{\otimes 2} \right]. \end{aligned}$$

Further, we have

$$\begin{aligned} & \frac{\partial \mathcal{L}(\Theta, \beta)}{\partial \beta} \\ = & -(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} \left[ \left\{ \frac{f_2(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})}{f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})} \mathbf{X}_{ij} - \frac{\int f_2(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) \mathbf{X} d\mathbf{X}}{\int f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X}} \right\} \right], \\ & \frac{\partial^2 \mathcal{L}(\Theta, \beta)}{\partial \beta \partial \beta^T} \\ = & -(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} \left[ \left\{ \frac{f_{22}(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})}{f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})} - \frac{f_2^2(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})}{f^2(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})} \right\} \mathbf{X}_{ij}^{\otimes 2} \right. \\ & \left. - \left\{ \frac{\int f_{22}(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) \mathbf{X}^{\otimes 2} d\mathbf{X}}{\int f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X}} - \frac{\left\{ \int f_2(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) \mathbf{X} d\mathbf{X} \right\}^{\otimes 2}}{\left\{ \int f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X} \right\}^2} \right\} \right], \end{aligned}$$

and

$$\begin{aligned} & \frac{\partial^2 \mathcal{L}(\Theta, \beta)}{\partial \text{vec}(\Theta) \partial \beta^T} \\ = & -(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} \left[ \left\{ \frac{f_{22}(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})}{f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})} - \frac{f_2^2(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})}{f^2(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}_{ij})} \right\} \text{vec}(\mathbf{z}_{ij}) \mathbf{X}_{ij}^T \right. \\ & \left. - \left\{ \frac{\int f_{22}(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) \text{vec}(\mathbf{z}_{ij}) \mathbf{X}^T d\mathbf{X}}{\int f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X}} \right. \right. \\ & \left. \left. - \frac{\left\{ \int f_2(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X} \text{vec}(\mathbf{z}_{ij}) \right\} \left\{ \int f_2(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) \mathbf{X} d\mathbf{X} \right\}^T}{\left\{ \int f(Y_{ij}, \Theta_{ij} + \beta^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X} \right\}^2} \right\} \right]. \end{aligned}$$



Let  $\mathbf{M} = (\text{vec}(\boldsymbol{\Theta})^\top, \boldsymbol{\beta}^\top)^\top$ , we can write

$$\frac{\partial \mathcal{L}(\boldsymbol{\Theta}, \boldsymbol{\beta})}{\partial \mathbf{M}} = -(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} [\mathbf{S}(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}, \boldsymbol{\beta}) - E\{\mathbf{S}(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}, \boldsymbol{\beta}) | Y_{ij}\}].$$

Furthermore because

$$\mathbf{H}(Y_{ij}, \mathbf{X}_{ij}, \boldsymbol{\Theta}, \boldsymbol{\beta}) = - \left\{ \frac{f_{22}(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}_{ij})}{f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}_{ij})} - \frac{f_2^2(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}_{ij})}{f^2(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}_{ij})} \right\} [\{\text{vec}(\mathbf{z}_{ij})^\top, \mathbf{X}_{ij}^\top\}^\top]^\otimes 2,$$

we get

$$\begin{aligned} E\{\mathbf{H}(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}, \boldsymbol{\beta}) | Y_{ij}\} &= - \frac{\int f_{22}(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}) g(\mathbf{X}) \{(\text{vec}(\mathbf{z}_{ij})^\top, \mathbf{X}^\top)^\top\}^\otimes 2 d\mathbf{X}}{\int f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}) g(\mathbf{X}) d\mathbf{X}} \\ &\quad + E \left\{ \frac{f_2^2(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X})}{f^2(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X})} \{(\text{vec}(\mathbf{z}_{ij})^\top, \mathbf{X}^\top)^\top\}^\otimes 2 | Y_{ij} \right\}, \end{aligned}$$

and

$$\begin{aligned} \frac{\partial^2 \mathcal{L}(\boldsymbol{\Theta}, \boldsymbol{\beta})}{\partial \mathbf{M} \partial \mathbf{M}^\top} &= (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} \left( \mathbf{H}(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}, \boldsymbol{\beta}) - E\{\mathbf{H}(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}, \boldsymbol{\beta}) | Y_{ij}\} \right. \\ &\quad \left. + E\{\mathbf{S}(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}, \boldsymbol{\beta})^\otimes 2 | Y_{ij}\} - [E\{\mathbf{S}(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}, \boldsymbol{\beta}) | Y_{ij}\}]^\otimes 2 \right). \end{aligned} \quad (5)$$

Now for any functions  $h_1(R_{ij}), h_2(Y_{ij}, \mathbf{X}_{ij})$ , because  $\mathbf{X}_{ij}$  and  $R_{ij}$  are independent given  $Y_{ij}$ , we have

$$E\{h_1(R_{ij}) h_2(Y_{ij}, \mathbf{X}_{ij}) | Y_{ij}\} = E\{h_1(R_{ij}) | Y_{ij}\} E\{h_2(Y_{ij}, \mathbf{X}_{ij}) | Y_{ij}\}.$$

We also have

$$E(R_{ij} [\mathbf{H}(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}, \boldsymbol{\beta}) - E\{\mathbf{H}(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}, \boldsymbol{\beta}) | Y_{ij}\}]) = \mathbf{0}.$$

Combining with (5), we obtain that  $E\{\partial^2 \mathcal{L}(\boldsymbol{\Theta}, \boldsymbol{\alpha}) / \partial \mathbf{M} \partial \mathbf{M}^\top\}$  is semi-positive definite. That is

$$\alpha_{\min}[E\{\partial^2 \mathcal{L}(\boldsymbol{\Theta}, \boldsymbol{\alpha}) / \partial \mathbf{M} \partial \mathbf{M}^\top\}] \geq 0, \quad (6)$$

where  $\alpha_{\min}(\mathbf{A})$  is the smallest singular value of matrix  $\mathbf{A}$ . In addition, we write

$$\begin{aligned} W_{ijk}(\boldsymbol{\Theta}, \boldsymbol{\beta}) &\equiv \mathbf{e}_k^\top \frac{\partial \ell_{ij}(\boldsymbol{\Theta}, \boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial (\boldsymbol{\Theta}, \mathbf{z}_{ij})} \\ &= \left[ \left\{ \frac{f_{22}(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}_{ij})}{f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}_{ij})} - \frac{f_2^2(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}_{ij})}{f^2(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}_{ij})} \right\} X_{ijk} \right. \\ &\quad - \left\{ \frac{\int f_{22}(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}) g(\mathbf{X}) X_k d\mathbf{X}}{\int f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}) g(\mathbf{X}) d\mathbf{X}} \right. \\ &\quad \left. \left. - \frac{\{\int f_2(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}) g(\mathbf{X}) X_k d\mathbf{X}\} \{\int f_2(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}) g(\mathbf{X}) d\mathbf{X}\}}{\left\{ \int f(Y_{ij}, \boldsymbol{\Theta}_{ij} + \boldsymbol{\beta}^\top \mathbf{X}) g(\mathbf{X}) d\mathbf{X} \right\}^2} \right\} \right]. \end{aligned}$$

## Appendix B. Some Lemmas

**Lemma A.1** *Suppose the independent random variables  $X_i, i = 1, \dots, n$  are bounded in the interval  $[a, b]$  with mean  $\mu$ . Then  $X_i$  is a sub-Gaussian with parameter  $\sigma = |b - a|$ . Further,*

$$\Pr \left\{ \left| \sum_{i=1}^n (X_i - \mu) \right| \geq t \right\} \leq 2 \exp \left\{ -\frac{t^2}{2n\sigma^2} \right\}.$$

Proof: This is the direct consequence of Example 2.3 and Proposition 2.1 in Chapter 2 of Wainwright (2019).

The lemma shows that a bounded variable is a sub-Gaussian distributed random variable.

**Lemma A.2** *Consider the independent random variables  $Y_1, \dots, Y_n$  such that there are infinitely many constants  $u_i, v_i, u_i \leq Y_i \leq v_i, i = 1, \dots, n$ . Let  $Z = \sup_{\mathbf{t} \in \mathcal{T}} \sum_{i=1}^n t_i Y_i$ , where  $\mathcal{T}$  is a set of vectors  $\mathbf{t} = (t_1, \dots, t_n)$  and  $\sigma = \sup_{\mathbf{t} \in \mathcal{T}} \left\{ \sum_{i=1}^n t_i^2 (v_i - u_i)^2 \right\}^{1/2} < \infty$ . Let  $m_Z$  be the median of  $Z$ . Then for  $\delta \geq 0$ , we have*

$$\Pr(|Z - m_Z| \geq \delta) \leq 4 \exp\{-\delta^2/(4\sigma^2)\}. \quad (7)$$

Furthermore

$$\begin{aligned} |E(Z) - m_Z| &\leq 4\sqrt{\pi}\sigma \\ \text{var}(Z) &\leq 16\sigma^2. \end{aligned} \quad (8)$$

Proof: This lemma is the Corollary 4.8 in Ledoux (2001).

The lemma demonstrates that a random variable with bounded variation has similar tail property as a sub-Gaussian distributed random variable does. It also provides bounds on the deviation from median and the distance between mean and median.

**Lemma A.3** *Let  $\mathbf{W}^i$  be independent  $d_r \times d_c$  zero mean random matrix such that  $\|\mathbf{W}^i\|_{op} \leq M$ , and define*

$$\sigma_i^2 \equiv \max \left\{ \|E(\mathbf{W}^i \mathbf{W}^{iT})\|_{op}, \|E(\mathbf{W}^{iT} \mathbf{W}^i)\|_{op} \right\}$$

and  $\sigma^2 = \sum_{i=1}^n \sigma_i^2$ . Then we have

$$\Pr \left( \left\| \sum_{i=1}^n \mathbf{W}^i \right\|_{op} \geq t \right) \leq d_r d_c \max \left[ \exp\{-t^2/(4\sigma^2)\}, \exp\{-t/(2M)\} \right].$$

Proof: This lemma follows Lemma 7 in Negahban and Wainwright (2012).

**Lemma A.4 (Hoeffding bound).** *Let  $X_1, \dots, X_N$  be independent centered sub-Gaussian random variables, let  $K = \max_i \|X_i\|_{\psi_2}$ , where  $\|X_i\|_{\psi_2} \equiv \sup_{k \geq 1} k^{-1/2} E(|X|^k)^{1/k}$ . is bounded for the sub-Gaussian random variable. Then for  $t > 0$ , we have*

$$\Pr \left\{ \left| \sum_{i=1}^N X_i \right| > t \right\} \leq 2 \exp \left( -\frac{t^2}{2NK^2} \right).$$

Proof: This lemma is from Proposition 2.1 in Wainwright (2019).

The lemma demonstrates the tail property for the operation norm of the sum of random matrices.

## Appendix C. Lemmas for Theorem 1

Define

$$\begin{aligned}\mathcal{C}(\gamma) &= \{\boldsymbol{\beta} \in \mathbf{R}^p : \|\boldsymbol{\beta}\|_\infty \|\boldsymbol{\beta}\|_1 / \|\boldsymbol{\beta}\|_2^2 \leq \gamma^{-1} \{mn / \log\{\max(p, mn)\}\}^{1/2}\}, \\ \mathbb{B}(D) &= \{\boldsymbol{\Delta}_\beta \in \mathcal{C}(\gamma) : \|\boldsymbol{\Delta}_\beta\|_2 \leq D, \|\boldsymbol{\Delta}_\beta\|_1 \leq D^2 / \gamma \sqrt{mn / \log\{\max(p, mn)\}}, \|\boldsymbol{\Delta}_\beta\|_\infty \leq a\}.\end{aligned}$$

**Lemma A.5** *Assume Condition (C1) and (C3) hold, and let  $\widehat{\boldsymbol{\Theta}}$  be the solution for (1), then*

$$\left\| \frac{\partial \mathcal{L}(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\beta}} - \frac{\partial \mathcal{L}(\boldsymbol{\Theta}_0, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\beta}} \right\|_\infty \leq 2ad_{\mathbf{W}}(mn)^{-1}$$

Proof: First note that by the mean value theorem, we have

$$\frac{\partial \mathcal{L}(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\beta}} - \frac{\partial \mathcal{L}(\boldsymbol{\Theta}_0, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\beta}} = (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m \frac{\partial^2 \mathcal{L}_{ij}(\boldsymbol{\Theta}^*, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\beta}_0 \partial \langle \boldsymbol{\Theta}, \mathbf{z}_{ij} \rangle} \langle \mathbf{z}_{ij}, \widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}_0 \rangle,$$

where  $\boldsymbol{\Theta}^*$  is the point on the line connecting  $\widehat{\boldsymbol{\Theta}}$  and  $\boldsymbol{\Theta}_0$ . Hence

$$\begin{aligned}& \left\| \frac{\partial \mathcal{L}(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\beta}} - \frac{\partial \mathcal{L}(\boldsymbol{\Theta}_0, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\beta}} \right\|_\infty \\ &= \sup_k |(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} \langle W_{ijk}(\boldsymbol{\Theta}^*, \boldsymbol{\beta}_0) \mathbf{z}_{ij}, \widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}_0 \rangle| \\ &= \sup_k |(mn)^{-1} \langle \mathbf{R} \circ \mathbf{W}_k(\boldsymbol{\Theta}^*, \boldsymbol{\beta}_0), \widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}_0 \rangle| \\ &\leq (mn)^{-1} \sup_k \|\text{vec}\{\mathbf{R} \circ \mathbf{W}_k(\boldsymbol{\Theta}^*, \boldsymbol{\beta}_0)\}\|_1 \|\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}_0\|_{\max} \\ &\leq 2ad_{\mathbf{W}}(mn)^{-1}.\end{aligned}$$

The last inequality holds by Condition (C1) and (C3). This proves the result.

**Lemma A.6** *Assume Conditions (C1), (C4) and (C5) hold, there is a constant  $\omega > 0$  such that*

$$\left\| \frac{\partial \mathcal{L}(\boldsymbol{\Theta}_0, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\beta}} \right\|_\infty \leq \sqrt{\omega \log\{\max(p, mn)\} / (mn)} + (mn)^{-1} d_{EX}$$

with probability at least  $1 - 2\{\max(p, mn)\}^{-1}$ .

Proof: Recall that  $\mathbf{X}^k$  is the  $n \times m$  matrix with the  $(i, j)$ th element  $X_{ijk}$ , and  $S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}, \boldsymbol{\beta})$  is the partial derivative of  $-\log\{f(Y_{ij}, \boldsymbol{\Theta}_{0ij} + \boldsymbol{\beta}^T \mathbf{X}_{ij})\}$  with respect to  $\boldsymbol{\Theta}_{0ij} + \boldsymbol{\beta}^T \mathbf{X}_{ij}$ . Let  $\mathbf{S}_2(\mathbf{Y}, \mathbf{X} | \boldsymbol{\Theta}, \boldsymbol{\beta})$  be the  $n \times m$  matrix with the  $(i, j)$ th element  $S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}, \boldsymbol{\beta})$ , we write

$$\begin{aligned}& \sup_k \left| \left\{ \frac{\partial \mathcal{L}(\boldsymbol{\Theta}_0, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\beta}} \right\}_k \right| \\ &= \sup_k |(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} \left[ \left\{ \frac{f_2(Y_{ij}, \boldsymbol{\Theta}_{0ij} + \boldsymbol{\beta}_0^T \mathbf{X}_{ij})}{f(Y_{ij}, \boldsymbol{\Theta}_{0ij} + \boldsymbol{\beta}_0^T \mathbf{X}_{ij})} X_{ijk} - \frac{\int f_2(Y_{ij}, \boldsymbol{\Theta}_{0ij} + \boldsymbol{\beta}_0^T \mathbf{X}) g(\mathbf{X}) X_k d\mathbf{X}}{\int f(Y_{ij}, \boldsymbol{\Theta}_{0ij} + \boldsymbol{\beta}_0^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X}} \right\} \right]| \\ &= \sup_k |(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} S_2(Y_{ij}, \mathbf{X}_{ij}, \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) X_{ijk} - R_{ij} E\{S_2(Y_{ij}, \mathbf{X} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) X_k | Y_{ij}\})|\end{aligned}$$

$$\begin{aligned}
 &\leq \sup_k |(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) X_{ijk}| \\
 &\quad + \sup_k |(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) X_k | Y_{ij}\}| \\
 &\leq \sup_k |(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) X_{ijk}| + (mn)^{-1} d_{E\mathbf{X}}
 \end{aligned}$$

Here we used Condition (C4) in step 4 and the last step. Now because  $S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0)$  is sub-Gaussian and  $X_{ijk}$  are bounded as assumed in Conditions (C5),  $E_{sij} = R_{ij} S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) X_{ijk}$  is a sub-Gaussian random variable, by Lemma A.4, we have

$$\Pr \left( |(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E_{sij}| > t \right) \leq 2 \exp(-2mnt^2/\omega)$$

for some constant  $\omega > 0$ . Therefore, we have

$$\begin{aligned}
 &\Pr \left[ \sup_k |(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) X_{ijk}| \geq t \right] \\
 &\leq p \Pr(|(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E_{sij}| > t) \\
 &\leq 2 \exp\{-2mnt^2/\omega + \log(p)\}.
 \end{aligned}$$

Let  $t = \sqrt{\omega \log\{\max(p, mn)\}/(mn)}$ , we get

$$\Pr \left[ \sup_k |(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) X_{ijk}| \geq \sqrt{\omega \log\{\max(p, mn)\}/(mn)} \right] \leq 2\{\max(p, mn)\}^{-1}.$$

Plug in the above result to (9), we get

$$\left\| \frac{\partial \mathcal{L}(\boldsymbol{\Theta}_0, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\beta}} \right\|_{\infty} \geq \sqrt{\omega \log\{\max(p, mn)\}/(mn)} + (mn)^{-1} d_{E\mathbf{X}}$$

with probability at most  $2\{\max(p, mn)\}^{-1}$ .

**Lemma A.7** *Let  $N(\delta)$  be the  $\delta$ -covering number of  $\mathbb{B}(D)$ . Then there is a constant  $c_1 > 0$  such that*

$$\log\{N(\delta)\} \leq \frac{9}{\delta^2} c_1^2 D^4 mn / \gamma^2.$$

Here,  $\delta$ -covering number is defined as the number of disks with radius  $\delta$  and center in  $\mathbb{B}(D)$  needed to cover  $\mathbb{B}(D)$ .

Proof: Define  $\mathbb{B}_1(D) = \{\boldsymbol{\Delta}_{\beta} \in \mathbb{R}^p \mid \|\boldsymbol{\Delta}_{\beta}\|_1 \leq D^2 \sqrt{mn/\log\{\max(p, mn)\}}/\gamma\}$  and let  $\tilde{N}(\delta)$  be the  $\delta$ -covering number of  $\mathbb{B}_1(D)$ . We have  $\mathbb{B}(D) \subseteq \mathbb{B}_1(D)$  and  $N(\delta) \leq \tilde{N}(\delta)$ . Now by the Sudakov minoration (Theorem 5.6 in Pisier (1999)) for a  $p$ -dimensional vector  $\mathbf{G}$  containing independent identically distributed standard normal random variables  $G_i, i = 1, \dots, p$

$$\sqrt{\log\{\tilde{N}(\delta)\}} \leq \frac{3}{\delta} E \left( \sup_{\|\boldsymbol{\Delta}_{\beta}\|_1 \leq D^2 \sqrt{mn/\log\{\max(p, mn)\}}/\gamma} \langle \mathbf{G}, \boldsymbol{\Delta}_{\beta} \rangle \right)$$

$$\begin{aligned}
 &\leq \frac{3}{\delta} E(\max_{i \in [1, p]} |G_i|) D^2 \sqrt{mn/\log(p)}/\gamma \\
 &\leq \frac{3}{\delta} c_1 D^2 \sqrt{\log(p)} \sqrt{mn/\log(p)}/\gamma \\
 &= \frac{3}{\delta} c_1 D^2 \sqrt{mn}/\gamma,
 \end{aligned}$$

for some constant  $c_1 > 0$ . The second line holds by the duality of  $L_1$  and  $L_\infty$  norm. The third line holds because  $E(\sup_{i \in [1, n]} |Z_i|) = O_p\{\sqrt{\log(n)}\}$  for an independent identically distributed sequence of normal random variables  $Z_1, \dots, Z_n$ . Therefore, we have

$$\log\{N(\delta)\} \leq \frac{9}{\delta^2} c_1^2 D^4 mn/\gamma^2.$$

In the following we first show for  $\tilde{\Theta}$  and  $\tilde{\beta}$  in the feasible set that  $\|\tilde{\Theta}\|_{\max} \leq a$  and  $\|\tilde{\beta}\|_\infty \leq a$ , we have

$$\begin{aligned}
 &\sup_{\beta \in \mathbb{B}(D)} |\beta^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \beta - \beta^\top E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \\
 &\quad - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}] \beta|
 \end{aligned}$$

is upper bounded by a  $O(1)$  constants. Then we use the peeling arguments (Negahban and Wainwright, 2012) to show the boundedness of

$$\begin{aligned}
 &\sup_{\beta \in \mathcal{C}(\gamma)} |\beta^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \beta - \beta^\top E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \\
 &\quad - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}] \beta|.
 \end{aligned}$$

Lemma A.8 and Lemma A.9 are the auxiliary results to establish the first relation.

**Lemma A.8** *Assume Conditions (C1), (C5), (C6). Let  $\delta = D/\xi$ , where  $\xi$  is an absolute constant, and  $c_1$  is defined in Lemma A.7. Let  $\beta \in \mathbb{B}(D)$ , and  $\beta^k, k = 1, \dots, N(\delta)$  be a  $\delta$ -covering of  $\mathbb{B}(D)$  in  $L_2$  norm. By definition, given an arbitrary  $\beta \in \mathbb{B}(D)$ , there is some index  $k \in \{1, \dots, N(\delta)\}$  and a difference matrix  $\Delta_\beta$  with  $\|\Delta_\beta\|_2 \leq \delta$  such that  $\beta = \beta^k + \Delta_\beta$ . Recall that for  $\tilde{\Theta}$  and  $\tilde{\beta}$  in the feasible set that  $\|\tilde{\Theta}\|_{\max} \leq a$  and  $\|\tilde{\beta}\|_\infty \leq a$*

$$\begin{aligned}
 &\mathbf{a}^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \mathbf{b} \\
 &= (mn)^{-1} \left( \sum_{i=1}^n \sum_{j=1}^m \mathbf{a}^\top R_{ij} [H_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top - E\{H_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \right. \\
 &\quad \left. + E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} - E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}] \mathbf{b} \right).
 \end{aligned}$$

Then there are positive constants  $c_2, c_3$  such that

$$\begin{aligned}
 &\Pr \left( \sup_{k=1, \dots, N(\delta)} |\beta^{k^\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \beta^k - \beta^{k^\top} E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \right. \\
 &\quad \left. - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}] \beta^k| \right. \\
 &\quad \left. \geq \delta^2 + \frac{32(d_{H_2} + d_{S_2}^2) a^2 c_0^2}{\sqrt{mn}} \right) \\
 &\leq 4 \exp(-c_2 D^2 mn),
 \end{aligned}$$

$$\begin{aligned}
 & \Pr \left( \sup_{k=1, \dots, N(\delta)} |\boldsymbol{\beta}^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}})(\boldsymbol{\beta} - \boldsymbol{\beta}^k) - \boldsymbol{\beta}^{k\top} E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\}] \right. \\
 & \quad \left. - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}\right] (\boldsymbol{\beta} - \boldsymbol{\beta}^k)| \\
 & \quad \geq \delta^2 + \frac{64(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}} \Bigg) \\
 & \leq 4 \exp(-c_3 D^2 mn).
 \end{aligned}$$

Proof: First we write

$$\begin{aligned}
 & |\boldsymbol{\beta}^{k\top} R_{ij}[H_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top - E\{H_2(Y_{ij}, \mathbf{X} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \mathbf{X} \mathbf{X}^\top | Y_{ij}\}] \\
 & \quad + E\{S_2^2(Y_{ij}, \mathbf{X} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \mathbf{X} \mathbf{X}^\top | Y_{ij}\} - E\{S_2(Y_{ij}, \mathbf{X} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \mathbf{X} | Y_{ij}\}^{\otimes 2}] \boldsymbol{\beta}^k| \\
 & \leq \sup_{\mathbf{X}_{ij}} \{2d_{H_2} \|\boldsymbol{\beta}^k\|_\infty \|\mathbf{X}_{ij}\|_1 \|\boldsymbol{\beta}^k\|_\infty \|\mathbf{X}_{ij}\|_1 + 2d_{S_2}^2 \|\boldsymbol{\beta}^k\|_\infty \|\mathbf{X}_{ij}\|_1 \|\boldsymbol{\beta}^k\|_\infty \|\mathbf{X}_{ij}\|_1\} \\
 & = 2(d_{H_2} + d_{S_2}^2)(ac_0)^2.
 \end{aligned}$$

Also we have

$$\begin{aligned}
 & |\boldsymbol{\beta}^{k\top} R_{ij}[H_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top - E\{H_2(Y_{ij}, \mathbf{X} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \mathbf{X} \mathbf{X}^\top | Y_{ij}\}] \\
 & \quad + E\{S_2^2(Y_{ij}, \mathbf{X} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \mathbf{X} \mathbf{X}^\top | Y_{ij}\} - E\{S_2(Y_{ij}, \mathbf{X} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \mathbf{X} | Y_{ij}\}^{\otimes 2}] \boldsymbol{\beta}^k| \\
 & \leq \sup_{\mathbf{X}_{ij}} \{2d_{H_2} \|\boldsymbol{\beta}^k\|_2 \|\mathbf{X}_{ij}\|_2 \|\boldsymbol{\beta}^k\|_\infty \|\mathbf{X}_{ij}\|_1 + 2d_{S_2}^2 \|\boldsymbol{\beta}^k\|_2 \|\mathbf{X}_{ij}\|_2 \|\boldsymbol{\beta}^k\|_\infty \|\mathbf{X}_{ij}\|_1\} \\
 & = 2(d_{H_2} + d_{S_2}^2)(ac_0^2)D,
 \end{aligned}$$

by Conditions (C1), (C5) and (C6). We now use Lemma A.2, where we treat each summand in  $\boldsymbol{\beta}^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \boldsymbol{\beta}^k$  as  $Y_i$  in Lemma A.2, and set  $(u_i, v_i) = \{-2(d_{H_2} + d_{S_2}^2)ac_0^2D, 2(d_{H_2} + d_{S_2}^2)ac_0^2D\}$  and  $(u_{1i}, v_{1i}) = \{-2(d_{H_2} + d_{S_2}^2)a^2c_0^2, 2(d_{H_2} + d_{S_2}^2)a^2c_0^2\}$ . Consider  $\mathcal{T}$  to contain only one vector  $\mathbf{t}$ , where  $\mathbf{t}$  is the  $mn$  dimensional vector with element  $(mn)^{-1}$ . Then  $Z$  in Lemma A.2 is  $\boldsymbol{\beta}^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \boldsymbol{\beta}^k$  and we set  $\sigma$  in Lemma A.2 (7) as  $\sigma = [\sum_{i=1}^n \sum_{j=1}^m (mn)^{-2} \{4(d_{H_2} + d_{S_2}^2)ac_0^2D\}^2]^{1/2} = 4(d_{H_2} + d_{S_2}^2)ac_0^2D/\sqrt{mn}$ , while we set  $\sigma$  in (8) as  $\sigma_1 = [\sum_{i=1}^n \sum_{j=1}^m (mn)^{-2} \{4(d_{H_2} + d_{S_2}^2)a^2c_0^2\}^2]^{1/2} = 4(d_{H_2} + d_{S_2}^2)a^2c_0^2/\sqrt{mn}$ . Let  $m_F$  be the median of  $\boldsymbol{\beta}^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \boldsymbol{\beta}^k$ , we have

$$\begin{aligned}
 & \Pr \left( |\boldsymbol{\beta}^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \boldsymbol{\beta}^k - \boldsymbol{\beta}^{k\top} E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}})\} \boldsymbol{\beta}^k| \geq \delta^2 + \frac{32(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}} \right) \\
 & \leq \Pr \left( |\boldsymbol{\beta}^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \boldsymbol{\beta}^k - m_F| + |\boldsymbol{\beta}^{k\top} E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}})\} \boldsymbol{\beta}^k - m_F| \geq \delta^2 + \frac{32(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}} \right) \\
 & \leq \Pr \left( |\boldsymbol{\beta}^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \boldsymbol{\beta}^k - m_F| \geq \delta^2 + \frac{32(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}} - \frac{16\sqrt{\pi}(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}} \right) \\
 & \leq \Pr \left( |\boldsymbol{\beta}^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \boldsymbol{\beta}^k - m_F| \geq \delta^2 \right) \\
 & \leq 4 \exp \left\{ \frac{-mn\delta^4}{64(d_{H_2} + d_{S_2}^2)^2 a^2 c_0^4 D^2} \right\}.
 \end{aligned}$$

Now for  $\boldsymbol{\beta}^k \in \mathbb{B}(D)$ ,  $k = 1, \dots, N(\delta)$ , we have

$$\Pr \left\{ \sup_{k=1, \dots, N(\delta)} |\boldsymbol{\beta}^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \boldsymbol{\beta}^k - \boldsymbol{\beta}^{k\top} E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\beta}}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\}] \right.$$

$$\begin{aligned}
 & -R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2}\beta^k \geq \delta^2 + \frac{32(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}} \Big\} \\
 \leq & 4 \exp \left\{ \frac{-mn\delta^4}{64(d_{H_2} + d_{S_2}^2)^2a^2c_0^4D^2} + \log\{N(\delta)\} \right\} \\
 \leq & 4 \exp \left\{ \frac{-mn\delta^4}{64(d_{H_2} + d_{S_2}^2)^2a^2c_0^4D^2} + \frac{9}{\delta^2}c_1^2D^4mn/\gamma^2 \right\}.
 \end{aligned}$$

Now because  $\delta = D/\xi$ , we can select  $\gamma$  sufficiently large so that  $\delta^4\{64(d_{H_2} + d_{S_2}^2)^2a^2c_0^4D^2\}^{-1} \geq 18\delta^{-2}c_1^2D^4/\gamma^2$ . Thus we have

$$\begin{aligned}
 & \Pr \left\{ \sup_{k=1, \dots, N(\delta)} |\beta^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\beta^k - \beta^{k\top} E\{R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}\mathbf{X}_{ij}^\top|Y_{ij}\} \right. \\
 & \quad \left. - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2}\beta^k \geq \delta^2 + \frac{32(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}} \right\} \\
 \leq & 4 \exp(-c_2D^2mn),
 \end{aligned}$$

where  $c_2 = -9c_1^2\xi^2\gamma^{-2}$ . Using the similar argument we can show that there is a  $c_3 > 0$  such that

$$\begin{aligned}
 & \Pr \left\{ \sup_{k=1, \dots, N(\delta)} |\beta^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y})(\beta - \beta^k) - \beta^{k\top} E\{R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}\mathbf{X}_{ij}^\top|Y_{ij}\} \right. \\
 & \quad \left. - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2}\beta^k\}(\beta - \beta^k)| \right. \\
 & \quad \left. \geq D^2/\xi^2 + \frac{64(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}} \right\} \\
 \leq & 4 \exp(-c_3D^2mn).
 \end{aligned}$$

This proves the result.

**Lemma A.9** *Assume Conditions (C5) and (C6) hold,  $\delta = D/\xi$ . Let  $\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})$  be as defined in Lemma A.8. Define  $\mathcal{D}_\beta(D) \equiv \{\Delta_\beta \in \mathbb{R}^p \mid \|\Delta_\beta\|_2 \leq \delta, \|\Delta_\beta\|_1 \leq 2D^2\sqrt{mn}/\log\{\max(p, mn)\}/\gamma, \|\Delta_\beta\|_\infty \leq 2a\}$ . Then*

$$\sup_{\Delta_\beta \in \mathcal{D}_\beta(D)} |\Delta_\beta^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\Delta_\beta| \leq 2(d_{H_2} + d_{S_2}^2)c_0^2\delta^2.$$

Proof: First note that

$$\sup_{\Delta_\beta \in \mathcal{D}_\beta(D)} (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m \Delta_\beta^\top H_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}\mathbf{X}_{ij}^\top \Delta_\beta \leq d_{H_2}c_0^2\delta^2.$$

and first inequality holds by Conditions (C5) and (C6).

Follow the same arguments we have

$$\begin{aligned}
 & \sup_{\Delta_\beta \in \mathcal{D}_\beta(D)} \left\{ (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m \Delta_\beta^\top E\{H_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}\mathbf{X}_{ij}^\top|Y_{ij}\} \Delta_\beta \right\} \leq d_{H_2}c_0^2\delta^2, \\
 & \sup_{\Delta_\beta \in \mathcal{D}_\beta(D)} \left\{ (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m \Delta_\beta^\top E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}\mathbf{X}_{ij}^\top|Y_{ij}\} \Delta_\beta \right\} \leq d_{S_2}^2c_0^2\delta^2, \tag{9}
 \end{aligned}$$

and

$$\sup_{\Delta_\beta \in \mathcal{D}_\beta(D)} \left\{ (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m \Delta_\beta^\top E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\}^{\otimes 2} \Delta_\beta \right\} \leq d_{S_2}^2 c_0^2 \delta^2 \quad (10)$$

Hence

$$\sup_{\Delta_\beta \in \mathcal{D}_\beta(D)} |\Delta_\beta^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \Delta_\beta| \leq 2(d_{H_2} + d_{S_2}^2) c_0^2 \delta^2.$$

**Lemma A.10** *Assume Conditions (C1)–(C6) hold. Let  $\delta$  be as defined in Lemma A.9. For  $\beta \in \mathcal{C}(\gamma)$ , we have*

$$\begin{aligned} & |\beta^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \beta - \beta^\top E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \\ & \quad - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}] \beta| \\ & \geq \alpha_{\min}(E[R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \\ & \quad - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}]) / 2 \|\beta\|_2^2 + 160(d_{H_2} + d_{S_2}^2) a^2 c_0^2 / \sqrt{mn} \end{aligned}$$

with probability at most  $1 - \exp[-C \log\{\max(p, mn)\}]$  for some positive constant  $C$ .

Proof: For any  $\beta \in \mathbb{B}(D)$ , there is a  $k \in \{1, \dots, N(\delta)\}$ , so that  $\tilde{\Delta}_\beta \equiv \beta - \beta^k$  satisfies  $\|\tilde{\Delta}_\beta\|_2 \leq \delta$ ,  $\|\tilde{\Delta}_\beta\|_1 \leq 2D^2 \sqrt{mn/\log\{\max(p, mn)\}}/\gamma$ , and  $\|\tilde{\Delta}_\beta\|_\infty \leq 2a$ , hence  $\tilde{\Delta}_\beta \in \mathcal{D}_\beta(D)$ , where  $\mathcal{D}_\beta(D)$  is defined in Lemma A.9. In addition,

$$\begin{aligned} & \beta^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \beta - \beta^\top E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \\ & \quad - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}] \beta \\ & = (\beta^k + \tilde{\Delta}_\beta)^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) (\beta^k + \tilde{\Delta}_\beta) - (\beta^k + \tilde{\Delta}_\beta)^\top E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \\ & \quad - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}] (\beta^k + \tilde{\Delta}_\beta) \\ & = \beta^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \beta^k - (\beta^k)^\top E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \\ & \quad - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}] \beta^k \\ & \quad + 2\beta^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \tilde{\Delta}_\beta - (\beta^k)^\top E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \\ & \quad - 2R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}] \tilde{\Delta}_\beta \\ & \quad + \tilde{\Delta}_\beta^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \tilde{\Delta}_\beta - \tilde{\Delta}_\beta^\top E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \\ & \quad - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}] \tilde{\Delta}_\beta. \end{aligned}$$

Hence by Lemmas A.8 and A.9,

$$\begin{aligned} & \sup_{\beta \in \mathbb{B}(D)} |\beta^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \beta - \beta^\top E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \\ & \quad - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}] \beta| \\ & \leq \sup_{k=1, \dots, N(\delta)} |\beta^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \beta^k - (\beta^k)^\top E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \\ & \quad - E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}] \beta^k| \\ & \quad + 2 \sup_{k=1, \dots, N(\delta)} |\beta^{k\top} \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \tilde{\Delta}_\beta - (\beta^k)^\top E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \\ & \quad - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}] \tilde{\Delta}_\beta| \end{aligned}$$



$$\begin{aligned}
 & -R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2}\tilde{\Delta}_\beta| \\
 & + \sup_{\tilde{\Delta}_\beta \in \mathcal{D}_\beta(D)} |\tilde{\Delta}_\beta^T \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\tilde{\Delta}_\beta| + \sup_{\tilde{\Delta}_\beta \in \mathcal{D}_\beta(D)} |\tilde{\Delta}_\beta^T E[E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}\mathbf{X}_{ij}^T|Y_{ij}\} \\
 & - E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2}]\tilde{\Delta}_\beta| \\
 \leq & 3\delta^2 + \frac{160(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}} + (2d_{H_2}c_0^2 + 2d_{S_2}^2c_0^2)\delta^2 + 2d_{S_2}^2c_0^2\delta^2 \\
 = & (3 + 2d_{H_2}c_0^2 + 4d_{S_2}^2c_0^2)D^2/\xi^2 + \frac{160(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}} \\
 = & D_5D^2/\xi^2 + \frac{160(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}} \tag{11}
 \end{aligned}$$

with probability at least  $1 - b_3\{\exp(-b_2D^2mn)\}$ , where  $b_2, b_3$  are constants not depending on  $D$ , and

$$D_5 = 3 + 2d_{H_2}c_0^2 + 4d_{S_2}^2c_0^2.$$

The second to the last equality holds by Lemma A.8, A.9 and (9), (10) in Lemma A.9. Now note for  $\beta \in \mathcal{C}(\gamma)$  with  $\|\beta\|_\infty = b$  for a constant  $b \leq a$ , we have

$$\|\beta\|_2^2 \geq b\gamma\|\beta\|_1 \sqrt{\frac{\log\{\max(p, mn)\}}{mn}} \geq b\gamma\|\beta\|_2 \sqrt{\frac{\log\{\max(p, mn)\}}{mn}},$$

which implies  $\|\beta\|_2 \geq b\gamma\sqrt{\log\{\max(p, mn)\}/(mn)}$ . Define  $\mu(b) = b^2\gamma^2\log\{\max(p, mn)\}/(mn)$ . Further let

$$\alpha_{0\beta}(\tilde{\Theta}, \tilde{\beta}) = \alpha_{\min}(E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}\mathbf{X}_{ij}^T|Y_{ij}\} - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2}],$$

$\alpha = \alpha_{0\beta}(\tilde{\Theta}, \tilde{\beta})\xi^2/D_5$  and  $\xi \geq \sqrt{D_5/\alpha_{0\beta}}$  so that  $\alpha > 1$ . Define

$$\begin{aligned}
 \mathcal{S}_l(b) &= \{\beta \in \mathcal{C}(\gamma) \mid \|\beta\|_\infty = b, \\
 & \sqrt{\alpha^{l-1}\mu(b)} \leq \|\beta\|_2 \leq \sqrt{\alpha^l\mu(b)}, \text{ and} \\
 & \|\beta\|_1 \leq b^{-1}\alpha^l\mu(b)\gamma^{-1}\sqrt{mn/\log\{\max(p, mn)\}}\}.
 \end{aligned}$$

Then  $\mathcal{S}_l(b) \subseteq \mathbb{B}(\sqrt{\alpha^l\mu(b)})$  and  $\{\beta : \|\beta\|_\infty = b\} \cap \mathcal{C}(\gamma) \subset \cup_{l=1}^\infty \mathcal{S}_l(b)$ .

For  $\beta \in \mathcal{S}_l(b)$ , we have

$$\begin{aligned}
 \alpha_{0\beta}(\tilde{\Theta}, \tilde{\beta})\|\beta\|_2^2 + \frac{160(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}} &\geq \alpha_{0\beta}(\tilde{\Theta}, \tilde{\beta})\alpha^{l-1}\mu(b) + \frac{160(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}} \\
 &= \alpha_{0\beta}\alpha^{-1}\alpha^l\mu(b) + \frac{160(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}} \\
 &= D_5\alpha^l\mu(b)/\xi^2 + \frac{160(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}}.
 \end{aligned}$$

Therefore,

$$\begin{aligned}
 & \Pr \left\{ \|\beta^T \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\beta - \beta^T E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}\mathbf{X}_{ij}^T|Y_{ij}\} \right. \\
 & \left. - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2}]\beta \right\} \\
 & \geq \alpha_{\min}(E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}\mathbf{X}_{ij}^T|Y_{ij}\}
 \end{aligned}$$

$$\begin{aligned}
 & -R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2})\|\beta\|_2^2 + \frac{160(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}}, \beta \in \mathcal{C}(\gamma) \Big\} \\
 \leq & \int_0^a \Pr \left\{ |\beta^T \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\beta - \beta^T E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}\mathbf{X}_{ij}^T|Y_{ij}\} \right. \\
 & \left. - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2}]\beta| \right. \\
 & \left. \geq \alpha_{\min}(E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}\mathbf{X}_{ij}^T|Y_{ij}\} \right. \\
 & \left. - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2})\|\beta\|_2^2 + \frac{160(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}}, \beta \in \mathcal{C}(\gamma), \|\beta\|_\infty = b \right\} db \\
 \leq & \int_0^a \sum_{l=1}^\infty \Pr \left\{ |\beta^T \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\beta - \beta^T E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}\mathbf{X}_{ij}^T|Y_{ij}\} \right. \\
 & \left. - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2}]\beta| \right. \\
 & \left. \geq \alpha_{\min}(E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}\mathbf{X}_{ij}^T|Y_{ij}\} \right. \\
 & \left. - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2})\|\beta\|_2^2 + \frac{160(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}}, \beta \in \mathcal{S}_l(b) \right\} db \\
 = & \sum_{l=1}^\infty \Pr \left\{ |\beta^T \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\beta - \beta^T E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}\mathbf{X}_{ij}^T|Y_{ij}\} \right. \\
 & \left. - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2}]\beta| \geq D_5\alpha^l\mu(b^*)/\xi^2 + \frac{160(d_{H_2} + d_{S_2}^2)a^2c_0^2}{\sqrt{mn}}, \beta \in \mathcal{S}_l(b^*) \right\} \\
 \leq & b_3 \sum_{l=1}^\infty [\exp(-b_2\alpha^l\mu(b^*)mn)] \\
 \leq & b_3 \sum_{l=1}^\infty [\exp(-b_2l\mu(b^*)\log(\alpha)mn)] \\
 \leq & b_3 \frac{\exp(-b_2\log(\alpha)\mu(b^*)mn)}{1 - \exp(-b_2\log(\alpha)\mu(b^*)mn)} \\
 \leq & \exp[-C\log\{\max(p, mn)\}],
 \end{aligned}$$

where  $b^*$  is a point on the line connecting 0 and  $a$ . The fourth line holds by (11) with  $D = \sqrt{\alpha^l\mu(b^*)}$ , the fifth line holds because  $\alpha^l \geq l\log(\alpha)$  for  $\alpha > 1$ . The last equality holds because  $\mu(b^*) = b^{*2}\gamma^2\{\max(p, mn)\}/(mn)$ . This proves the result.

**Lemma A.11** *Assume Conditions (C5) and (C6) hold, suppose  $\Delta_\beta = \hat{\beta} - \beta_0$ , let  $\tilde{\Theta} = \hat{\Theta}$  and  $\tilde{\beta} = \beta^*$  a point connecting  $\hat{\beta}$  and  $\beta_0$ , then there is a  $\sigma_{1F} > 0$  such that*

$$\Delta_\beta^T \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\hat{\Theta}, \beta^*)\Delta_\beta \geq \Delta_\beta^T E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\hat{\Theta}, \beta^*)\}\Delta_\beta - \sigma_{1F}\sqrt{2\log\{\max(p, mn)\}/(mn)},$$

with probability at least  $1 - 2\max(p, mn)^{-1}$ , where  $\sigma_{1F} = 32c_0^2a^2(d_{H_2} + d_{S_2}^2)$  as defined in Theorem 1 and  $E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\hat{\Theta}, \beta^*)\} = E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\Theta, \beta)\}_{\Theta=\hat{\Theta}, \beta=\beta^*}$ .

Proof: First note that for  $\Delta_\beta$  with  $\|\Delta_\beta\|_\infty \leq 2a$ , we have

$$\begin{aligned}
 & |(\Delta_\beta)^T R_{ij}H_2(Y_{ij}, \mathbf{X}_{ij}|\hat{\Theta}, \beta^*)\mathbf{X}_{ij}\mathbf{X}_{ij}^T)^T \Delta_\beta| \\
 & \leq \{ \|\Delta_\beta\|_{\max} \|\mathbf{X}_{ij}\|_1 \}^2 H_2(Y_{ij}, \mathbf{X}_{ij}|\hat{\Theta}, \beta^*) \leq 4c_0^2a^2d_{H_2},
 \end{aligned}$$

Similarly,

$$|\Delta_\beta^T R_{ij}E\{H_2(Y_{ij}, \mathbf{X}_{ij}|\hat{\Theta}, \beta^*)\mathbf{X}_{ij}\mathbf{X}_{ij}^T|Y_{ij}\}\Delta_\beta| \leq 4c_0^2a^2d_{H_2},$$

$$|(\Delta_\beta)^\top R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\widehat{\Theta}, \beta^*) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} (\Delta_\beta)| \leq 4c_0^2 a^2 d_{S_2}^2, \quad (12)$$

and

$$|(\Delta_\beta)^\top R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\widehat{\Theta}, \beta^*) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2} (\Delta_\beta)| \leq 4c_0^2 a^2 d_{S_2}^2. \quad (13)$$

Hence we have each summand in  $\Delta_\beta^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\widehat{\Theta}, \beta^*) \Delta_\beta - \Delta_\beta^\top E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\widehat{\Theta}, \beta^*)\} \Delta_\beta$  is in the range of  $\{-16(c_0^2 a^2 d_{H_2} + c_0^2 a^2 d_{S_2}^2), 16(c_0^2 a^2 d_{H_2} + c_0^2 a^2 d_{S_2}^2)\}$ , and in turn is sub-Gaussian with parameter  $\sigma_{1F}^2 \equiv \{32(c_0^2 a^2 d_{H_2} + c_0^2 a^2 d_{S_2}^2)\}^2$  by Lemma A.1. Here the expectation is taken over  $\mathbf{X}$  and  $\mathbf{Y}$ . Therefore, by Lemma A.1, we have

$$\begin{aligned} & \Pr \left\{ |\Delta_\beta^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\widetilde{\Theta}, \widetilde{\beta}) \Delta_\beta - \Delta_\beta^\top E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\widetilde{\Theta}, \widetilde{\beta})\} \Delta_\beta| > t \right\} \\ &= \Pr \left\{ |mn \Delta_\beta^\top F(\mathbf{X}, \mathbf{Y}|\widetilde{\Theta}, \widetilde{\beta}) \Delta_\beta - mn \Delta_\beta^\top E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\widetilde{\Theta}, \widetilde{\beta})\} \Delta_\beta| > mnt \right\} \\ &\leq 2 \exp \left\{ -\frac{(mn)^2 t^2}{2mn \sigma_{1F}^2} \right\} \\ &= 2 \exp \left( -\frac{mnt^2}{2\sigma_{1F}^2} \right). \end{aligned}$$

Let  $t = \sigma_{1F} \sqrt{2 \log\{\max(p, mn)\}/(mn)}$ , we obtain

$$\begin{aligned} & \Pr \left\{ |\Delta_\beta^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\widehat{\Theta}, \beta^*) \Delta_\beta - \Delta_\beta^\top E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\widehat{\Theta}, \beta^*)\} \Delta_\beta| > \sigma_{1F} \sqrt{2 \log\{\max(p, mn)\}/(mn)} \right\} \\ &\leq 2 \max(p, mn)^{-1}. \end{aligned}$$

Hence we have

$$\Delta_\beta^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\widehat{\Theta}, \beta^*) \Delta_\beta \geq \Delta_\beta^\top E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\widehat{\Theta}, \beta^*)\} \Delta_\beta - \sigma_{1F} \sqrt{2 \log\{\max(p, mn)\}/(mn)},$$

with probability at least  $1 - 2 \max(p, mn)^{-1}$ .

**Lemma A.12** *Assume Condition (C5) and (C6) hold, suppose  $\Delta_\beta = \widehat{\beta} - \beta_0$ . Let  $\lambda_\beta \geq 2 \|\partial \mathcal{L}(\widehat{\Theta}, \beta_0)/\partial \beta\|_{op}$  and  $\sigma_{1F}$  be as defined in Theorem 1. Then either*

$$\begin{aligned} \|\Delta_\beta\|_2 &\leq \left\{ \alpha_{\min}(E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\widehat{\Theta}, \beta^*) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \right. \\ &\quad \left. - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\widehat{\Theta}, \beta^*) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}]\right\}^{-1/2} \{2\sigma_{1F}^2 \log\{\max(p, mn)\}/(mn)\}^{1/4}, \end{aligned}$$

or  $\|\Delta_\beta\|_1 \leq 4\sqrt{s} \|\Delta_\beta\|_2$  with probability at least  $1 - 2\{\max(p, mn)\}^{-1}$ .

Proof. First consider  $\Delta_\beta^\top E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\widehat{\Theta}, \beta^*)\} \Delta_\beta \leq \sigma_{1F} \sqrt{2 \log\{\max(p, mn)\}/(mn)}$ . Because

$$\begin{aligned} E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\widehat{\Theta}, \beta^*)\} &= E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\widehat{\Theta}, \beta^*) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \\ &\quad - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\widehat{\Theta}, \beta^*) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}], \end{aligned}$$

we have

$$\begin{aligned} & \alpha_{\min}(E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\widehat{\Theta}, \beta^*) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \\ & \quad - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\widehat{\Theta}, \beta^*) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}]) \|\Delta_\beta\|_2^2 \leq \sigma_{1F} \sqrt{2 \log\{\max(p, mn)\}/(mn)}, \end{aligned}$$

and hence

$$\|\Delta_\beta\|_2 \leq \left\{ \alpha_{\min}(E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\widehat{\Theta}, \beta^*) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} \right.$$

$$-R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\widehat{\Theta}, \beta^*)\mathbf{X}_{ij}|Y_{ij}\}^{\otimes 2})\}^{-1/2}\sigma_{1F}^{1/2}\{2\log\{\max(p, mn)\}/(mn)\}^{1/4}.$$

Now consider  $\Delta_{\beta}^T E\{\mathbf{F}_{\beta}(\mathbf{X}, \mathbf{Y}|\widehat{\Theta}, \beta^*)\}\Delta_{\beta} > \sigma_{1F}\sqrt{2\log\{\max(p, mn)\}/(mn)}$ . We have

$$0 < \Delta_{\beta}^T \mathbf{F}_{\beta}(\mathbf{X}, \mathbf{Y}|\widehat{\Theta}, \beta^*)\Delta_{\beta}/2 = -\frac{\partial\mathcal{L}(\widehat{\Theta}, \beta_0)}{\partial\beta^T}\Delta_{\beta} + \lambda_{\beta}\|\beta_0\|_1 - \lambda_{\beta}\|\widehat{\beta}\|_1 \quad (14)$$

with probability at least  $1 - 2\max(p, mn)^{-1}$ . Let  $S$  be the set of indices that  $\beta_{0j} \neq 0$ , and  $\mathbf{v}_S$  be the sub-vector of  $\mathbf{v}$  with the elements  $j \in S$ . Then we have

$$\begin{aligned} \|\beta_0 + \Delta_{\beta}\|_1 + \|\Delta_{\beta S}\|_1 &\geq \|\beta_0 + \Delta_{\beta} - \Delta_{\beta S}\|_1 \\ &= \|\beta_0 + \Delta_{\beta S^c}\|_1 \\ &= \|\beta_{0S}\|_1 + \|\Delta_{\beta S^c}\|_1. \end{aligned}$$

Hence

$$\begin{aligned} \|\beta_0 + \Delta_{\beta}\|_1 - \|\beta_0\|_1 &\geq \{\|\beta_0\|_1 - \|\Delta_{\beta S}\|_1\} + \|\Delta_{\beta S^c}\|_1 - \|\beta_0\|_1 \\ &= \|\Delta_{\beta S^c}\|_1 - \|\Delta_{\beta S}\|_1. \end{aligned} \quad (15)$$

Combine with (14), we have

$$\begin{aligned} 0 &\leq -\frac{\partial\mathcal{L}(\widehat{\Theta}, \beta_0)}{\partial\beta^T}\Delta_{\beta} + \lambda_{\beta}(\|\Delta_{\beta S}\|_1 - \|\Delta_{\beta S^c}\|_1) \\ &\leq \left\|\frac{\partial\mathcal{L}(\widehat{\Theta}, \beta_0)}{\partial\beta^T}\right\|_{op}(\|\Delta_{\beta S}\|_1 + \|\Delta_{\beta S^c}\|_1) + \lambda_{\beta}(\|\Delta_{\beta S}\|_1 - \|\Delta_{\beta S^c}\|_1) \\ &\leq \lambda_{\beta}/2(\|\Delta_{\beta S}\|_1 + \|\Delta_{\beta S^c}\|_1) + \lambda_{\beta}(\|\Delta_{\beta S}\|_1 - \|\Delta_{\beta S^c}\|_1) \\ &= 3\lambda_{\beta}/2\|\Delta_{\beta S}\|_1 - \lambda_{\beta}/2\|\Delta_{\beta S^c}\|_1 \end{aligned}$$

which implies  $\|\Delta_{\beta S^c}\|_1 \leq 3\|\Delta_{\beta S}\|_1$  and in turn  $\|\Delta_{\beta}\|_1 \leq 4\|\Delta_{\beta S}\|_1 \leq 4\sqrt{s}\|\Delta_{\beta}\|_2$  with probability at least  $1 - 2p^{-1}$ . This proves the result.

## Appendix D. Proof of Theorem 1

First, we note that from  $\lambda_{\beta} \geq 2\sqrt{\omega\log\{\max(p, mn)\}/(mn)} + 2(mn)^{-1}d_{\mathbf{EX}} + 4ad_{\mathbf{W}}(mn)^{-1}$  in the theorem statement and Lemma A.5 and A.6, we obtain that

$$\left\|\frac{\partial\mathcal{L}(\widehat{\Theta}, \beta_0)}{\partial\beta^T}\right\|_{\infty} \leq \frac{\lambda_{\beta}}{2}$$

with probability at least  $1 - 2\{\max(p, mn)\}^{-1}$ . Further, because  $\|\mathbf{A}\|_{op} \leq \|\mathbf{A}\|_{\infty}$  for any matrix  $\mathbf{A}$ , we get

$$\left\|\frac{\partial\mathcal{L}(\widehat{\Theta}, \beta_0)}{\partial\beta^T}\right\|_{op} \leq \frac{\lambda_{\beta}}{2}.$$

When  $\Delta_{\beta}^T E\{\mathbf{F}_{\beta}(\mathbf{X}, \mathbf{Y})\}\Delta_{\beta} \leq \sigma_{1F}\sqrt{2\log\{\max(p, mn)\}/(mn)}$ , from the proof of Lemma A.12, we know

$$\|\Delta_{\beta}\|_2 \leq (4\alpha_{0\beta})^{-1/2}\{2\sigma_{1F}^2\log\{\max(p, mn)\}/(mn)\}^{1/4}.$$

We discuss two cases when  $\Delta_\beta^\top E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y})\}\Delta_\beta \geq \sigma_{1F} \sqrt{2\log\{\max(p, mn)\}/(mn)}$ . Case I:  $\Delta_\beta \notin \mathcal{C}(\gamma)$ . Then

$$\|\Delta_\beta\|_2^2 \leq \|\Delta_\beta\|_\infty \|\Delta_\beta\|_1 \gamma \sqrt{\frac{\log\{\max(p, mn)\}}{mn}} \leq 8a\sqrt{s} \|\Delta_\beta\|_2 \gamma \sqrt{\frac{\log\{\max(p, mn)\}}{mn}}$$

with probability at least  $1 - 2 \max(p, mn)^{-1}$ , which implies  $\|\Delta_\beta\|_2 \leq 8a\sqrt{s}\gamma \sqrt{\log\{\max(p, mn)\}/(mn)}$  with probability at least  $1 - 2 \max(p, mn)^{-1}$ .

Case II:  $\hat{\beta} - \beta_0 \in \mathcal{C}(\gamma)$ . Because  $\hat{\Theta}, \hat{\beta}$  is the minimizer of  $\mathcal{L}(\Theta, \beta) + \lambda_\Theta \|\Theta\|_* + \lambda_\beta \|\beta\|_1$ , we have

$$\mathcal{L}(\hat{\Theta}, \hat{\beta}) - \mathcal{L}(\hat{\Theta}, \beta_0) = \frac{\partial \mathcal{L}(\hat{\Theta}, \beta_0)}{\partial \beta^\top} \Delta_\beta + \frac{1}{2} (\hat{\beta} - \beta_0)^\top \frac{\partial \mathcal{L}(\hat{\Theta}, \beta^*)}{\partial \beta \partial \beta^\top} (\hat{\beta} - \beta_0)$$

and

$$\begin{aligned} \Delta_\beta^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \hat{\Theta}, \beta^*) \Delta_\beta / 2 &= (\hat{\beta} - \beta_0)^\top \frac{\partial^2 \mathcal{L}(\hat{\Theta}, \beta^*)}{\partial \beta \partial \beta^\top} (\hat{\beta} - \beta^*) / 2 \\ &= \mathcal{L}(\hat{\Theta}, \hat{\beta}) - \mathcal{L}(\hat{\Theta}, \beta_0) - \frac{\partial \mathcal{L}(\hat{\Theta}, \beta_0)}{\partial \beta^\top} \Delta_\beta \\ &\leq -\frac{\partial \mathcal{L}(\hat{\Theta}, \beta_0)}{\partial \beta^\top} \Delta_\beta + \lambda_\beta \|\beta_0\|_1 - \lambda_\beta \|\hat{\beta}\|_1 \\ &\leq \left\| \frac{\partial \mathcal{L}(\hat{\Theta}, \beta_0)}{\partial \beta^\top} \right\|_\infty \|\Delta_\beta\|_1 + \lambda_\beta \|\beta_0\|_1 - \lambda_\beta \|\hat{\beta}\|_1 \\ &\leq \lambda_\beta / 2 \|\Delta_\beta\|_1 + \lambda_\beta \|\beta_0 - \hat{\beta}\|_1 \\ &\leq 6\lambda_\beta \sqrt{s} \|\Delta_\beta\|_2 \end{aligned} \tag{16}$$

with probability at least  $1 - 2 \max(p, mn)^{-1} - 2(mn)^{-1}$ . The first inequality holds by the second order mean value theorem for  $\mathcal{L}(\hat{\Theta}, \beta)$  on  $\beta$ . The third inequality holds by the fact that  $\lambda_\beta \geq 2\sqrt{\omega \log\{\max(p, mn)\}/(mn)} + 2(mn)^{-1} d_{E\mathbf{X}} + 4ad_{\mathbf{W}}(mn)^{-1}$  in the theorem statement and Lemma A.5 and A.6, and triangular inequality. In the last inequality, we used Lemma A.12. Further, by Lemma A.10, with probability at least  $1 - \exp\{-C \log\{\max(p, mn)\}\}$ ,

$$\begin{aligned} &|\Delta_\beta^\top / 2 \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \hat{\Theta}, \beta^*) \Delta_\beta / 2 - \Delta_\beta^\top / 2 E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \hat{\Theta}, \beta^*)\} \Delta_\beta / 2| \\ &\leq 2\alpha_{0\beta} \|\Delta_\beta / 2\|_2^2 + 160(d_{H_2} + d_{S_2}^2) a^2 c_0^2 / \sqrt{mn}, \end{aligned}$$

so

$$\begin{aligned} &\Delta_\beta^\top / 2 \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \hat{\Theta}, \beta^*) \Delta_\beta / 2 - \Delta_\beta^\top / 2 E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \hat{\Theta}, \beta^*)\} \Delta_\beta / 2 \\ &\geq -2\alpha_{0\beta} \|\Delta_\beta / 2\|_2^2 - 160(d_{H_2} + d_{S_2}^2) a^2 c_0^2 / \sqrt{mn}. \end{aligned}$$

Thus,

$$\begin{aligned} &\Delta_\beta^\top / 2 \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \hat{\Theta}, \beta^*) \Delta_\beta / 2 \\ &\geq \alpha_{\min} (E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \hat{\Theta}, \beta^*) \mathbf{X}_{ij} \mathbf{X}_{ij}^\top | Y_{ij}\} - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \hat{\Theta}, \beta^*) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2}]) \|\Delta_\beta / 2\|_2^2 \\ &\quad - 2^{-1} \alpha_{0\beta} \|\Delta_\beta\|_2^2 - 160(d_{H_2} + d_{S_2}^2) a^2 c_0^2 / \sqrt{mn} \\ &= 2^{-1} \alpha_{0\beta} \|\Delta_\beta\|_2^2 - 160(d_{H_2} + d_{S_2}^2) a^2 c_0^2 / \sqrt{mn} \end{aligned}$$

with probability at least  $1 - \exp\{-C \log\{\max(p, mn)\}\}$ , which implies

$$\Delta_\beta^\top \mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \hat{\Theta}, \beta^*) \Delta_\beta / 2 \geq \alpha_{0\beta} \|\Delta_\beta\|_2^2 - 320(d_{H_2} + d_{S_2}^2) a^2 c_0^2 / \sqrt{mn},$$

with probability at least  $1 - \exp\{-C\log\{\max(p, mn)\}\}$ . Combining with (16), we have

$$\alpha_{0\beta}\|\Delta_{\beta}\|_2^2 \leq 6\lambda_{\beta}\sqrt{s}\|\Delta_{\beta}\|_2 + 320(d_{H_2} + d_{S_2}^2)a^2c_0^2/(\sqrt{mn})$$

with probability at least  $1 - 2\max(p, mn)^{-1} - 2(mn)^{-1} - \exp\{-C\log\{\max(p, mn)\}\}$ . Then

$$\|\Delta_{\beta}\|_2^2 - 6\lambda_{\beta}\sqrt{s}\|\Delta_{\beta}\|_2/\alpha_{0\beta} \leq 320(d_{H_2} + d_{S_2}^2)a^2c_0^2/(\alpha_{0\beta}\sqrt{mn}).$$

This leads to

$$(\|\Delta_{\beta}\|_2^2 - 3\lambda_{\beta}\sqrt{s}/\alpha_{0\beta})^2 \leq 320(d_{H_2} + d_{S_2}^2)a^2c_0^2/(\alpha_{0\beta}\sqrt{mn}) + (3\lambda_{\beta}\sqrt{s}/\alpha_{0\beta})^2,$$

Hence

$$\|\Delta_{\beta}\|_2 \leq \{320(d_{H_2} + d_{S_2}^2)a^2c_0^2/(\alpha_{0\beta}\sqrt{mn}) + (3\lambda_{\beta}\sqrt{s}/\alpha_{0\beta})^2\}^{1/2} + 3\lambda_{\beta}\sqrt{s}/\alpha_{0\beta},$$

with probability at least  $1 - 2\max(p, mn)^{-1} - 2(mn)^{-1} - \exp\{-C\log\{\max(p, mn)\}\}$ . Combine with the order in Case I and before Case I, we have

$$\begin{aligned} & \|\Delta_{\beta}\|_2 \\ & \leq \max\left(\{320(d_{H_2} + d_{S_2}^2)a^2c_0^2/(\alpha_{0\beta}\sqrt{mn}) + (3\lambda_{\beta}\sqrt{s}/\alpha_{0\beta})^2\}^{1/2} + 3\lambda_{\beta}\sqrt{s}/\alpha_{0\beta}, 8a\sqrt{s}\gamma\sqrt{\log\{\max(p, mn)\}/(mn)}, \right. \\ & \quad \left. (4\alpha_{0\beta})^{-1/2}\{2\sigma_{1F}^2\log\{\max(p, mn)\}/(mn)\}^{1/4}\right) \end{aligned}$$

with probability at least  $1 - 4\max(p, mn)^{-1} - 2(mn)^{-1} - 2\exp\{-C\log\{\max(p, mn)\}\}$ . This proves the result.

## Appendix E. Lemmas for Theorem 2

Let  $\rho(\nu, D) \equiv D^2/\{\nu\sqrt{d\log(d)/(mn)}\} = D^2/\{\nu\sqrt{\log(d)/d}\}$ . Define sets

$$\begin{aligned} \mathcal{C}_{\Theta}(\nu) & \equiv \left\{0 \neq \Theta \in \mathbb{R}^{m \times n} \mid \frac{\|\Theta\|_{\max}}{\|\Theta\|_F} \frac{\|\Theta\|_*}{\|\Theta\|_F} \leq \frac{1}{\nu} \sqrt{\frac{mn}{d\log(d)}}\right\} \\ & = \left\{0 \neq \Theta \in \mathbb{R}^{m \times n} \mid \frac{\|\Theta\|_{\max}}{\|\Theta\|_F} \frac{\|\Theta\|_*}{\|\Theta\|_F} \leq \frac{1}{\nu} \sqrt{\frac{d}{\log(d)}}\right\}, \\ \mathbb{B}_{\Theta}(D) & \equiv \{\Theta \in \mathcal{C}_{\Theta}(\nu) \mid \|\Theta\|_{\max} = a, \|\Theta\|_F \leq D, \|\Theta\|_* \leq \rho(\nu, D)\}, \\ \bar{\mathbb{B}}_{\Theta}(D) & \equiv \{\Theta \in \mathcal{C}_{\Theta}(\nu) \mid \|\Theta\|_{\max} \leq a, \|\Theta\|_F \leq D, \|\Theta\|_* \leq \rho(\nu, D)\}. \end{aligned}$$

Let  $\Theta^1, \dots, \Theta^{N_{\Theta}(\delta)}$  be a  $\delta$ -covering of  $\bar{\mathbb{B}}_{\Theta}(D)$  in Frobenius norm. By definition, given an arbitrary  $\Theta \in \bar{\mathbb{B}}_{\Theta}(D)$ , there is some index  $k \in \{1, \dots, N_{\Theta}(\delta)\}$  and a difference matrix  $\Delta_{\Theta}$  with  $\|\Delta_{\Theta}\|_F \leq \delta$  such that  $\Theta = \Theta^k + \Delta_{\Theta}$ .

**Lemma A.13** *Assume Conditions (C1) and (C6) hold. Then there is a positive constant  $c_d$  such that for sufficiently large  $n, m, d$ ,*

$$\Pr\left(\left\|\frac{\partial\mathcal{L}(\Theta_0, \beta_0)}{\partial\Theta}\right\|_{op} \geq c_d\sqrt{d\log(d)/(mn)}\right) < d^{-1},$$

or equivalently,

$$\Pr\left(\left\|\frac{\partial\mathcal{L}(\Theta_0, \beta_0)}{\partial\Theta}\right\|_{op} \geq c_d\sqrt{\log(d)/d}\right) < d^{-1}.$$

Proof: First note that

$$\begin{aligned}
 & \frac{\partial \mathcal{L}(\boldsymbol{\Theta}_0, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\Theta}} \\
 = & -(mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} \left[ \left\{ \frac{f_2(Y_{ij}, \boldsymbol{\Theta}_{0ij} + \boldsymbol{\beta}_0^T \mathbf{X}_{ij})}{f(Y_{ij}, \boldsymbol{\Theta}_{0ij} + \boldsymbol{\beta}_0^T \mathbf{X}_{ij})} - \frac{\int f_2(Y_{ij}, \boldsymbol{\Theta}_{0ij} + \boldsymbol{\beta}^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X}}{\int f(Y_{ij}, \boldsymbol{\Theta}_{0ij} + \boldsymbol{\beta}_0^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X}} \right\} \mathbf{e}_i \mathbf{e}_j^T \right] \\
 = & (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} [S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) - E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) | Y_{ij}\}] \mathbf{e}_i \mathbf{e}_j^T.
 \end{aligned}$$

Let  $\mathbf{W}^{ij} = R_{ij}[S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) - E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) | Y_{ij}\}] \mathbf{e}_i \mathbf{e}_j^T / (mn)$ . Then  $\mathbf{W}^{ij}$  is a mean zero random matrix with  $\|\mathbf{W}^{ij}\|_{op} \leq \sqrt{\|\mathbf{W}^{ij}\|_1 \|\mathbf{W}^{ij}\|_\infty} \leq 2d_{S_2} / (mn)$  by Condition (C6). Further

$$\begin{aligned}
 & \|E(\mathbf{W}^{ij} \mathbf{W}^{ijT})\|_{op} \\
 = & \|E(R_{ij}[S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) - E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}_0) | Y_{ij}\}]^2 \mathbf{e}_i \mathbf{e}_i^T / (mn)^2)\|_{op} \\
 \leq & \|4d_{S_2}^2 / (mn)^2 \mathbf{e}_i \mathbf{e}_i^T\|_{op} \\
 = & 4d_{S_2}^2 / (mn)^2.
 \end{aligned}$$

Similarly,

$$\|E(\mathbf{W}^{ijT} \mathbf{W}^{ij})\|_{op} \leq 4d_{S_2}^2 / (mn)^2.$$

Hence

$$\sigma_W^2 \equiv mn \max(\|E(\mathbf{W}^{ijT} \mathbf{W}^{ij})\|_{op}, \|E(\mathbf{W}^{ij} \mathbf{W}^{ijT})\|_{op}) \leq 4d_{S_2}^2 / (mn).$$

Therefore, by Lemma A.3 we have

$$\begin{aligned}
 & \Pr\left(\left\|\frac{\partial \mathcal{L}(\boldsymbol{\Theta}_0, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\Theta}}\right\|_{op} \geq t\right) \\
 \leq & \Pr\left(\left\|\sum_{i=1}^n \sum_{j=1}^m \mathbf{W}^{ij}\right\|_{op} \geq t\right) \\
 \leq & mn \max\left(\exp[-t^2 / \{16d_{S_2}^2 / (mn)\}], \exp[-t / \{4d_{S_2} / (mn)\}]\right).
 \end{aligned}$$

Now let  $t = c_d \sqrt{d \log(d) / (mn)}$  for  $c_d = 4d_{S_2}$ . Then we have

$$\begin{aligned}
 & \Pr\left(\left\|\frac{\partial \mathcal{L}(\boldsymbol{\Theta}_0, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\Theta}}\right\|_{op} \geq t\right) \\
 \leq & mn \max\left(\exp[-c_d^2 d \log(d) / (mn) / \{16d_{S_2}^2 / (mn)\}], \exp[-c_d \sqrt{d \log(d) / (mn)} / \{4d_{S_2} / (mn)\}]\right) \\
 \leq & mn \max\left(\exp\{-d \log(d)\}, \exp\{-mn \sqrt{d \log(d) / mn}\}\right) \\
 = & \exp\{-d \log(d) + \log(mn)\} \\
 \leq & d^{-1}
 \end{aligned}$$

for  $d, m, n \rightarrow \infty$ .

**Lemma A.14** *Assume Conditions (C1)–(C6) hold. Let*

$$\lambda_\beta \geq 2\sqrt{\omega \log\{\max(p, mn)\} / mn} + 2(mn)^{-1} d_{EX} + 4ad_{\mathbf{W}}(mn)^{-1}.$$

Then

$$\begin{aligned}
 & \left\| \frac{\partial \mathcal{L}(\boldsymbol{\Theta}_0, \widehat{\boldsymbol{\beta}})}{\partial \boldsymbol{\Theta}} - \frac{\partial \mathcal{L}(\boldsymbol{\Theta}_0, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\Theta}} \right\|_{op} \\
 \leq & (2d_{H_2} + 2d_{S_2}^2) \max \left( [320(d_{H_2} + d_{S_2}^2)a^2c_0^2/(\alpha_0\beta\sqrt{mn}) + \{3\lambda_\beta\sqrt{s}/\alpha_0\beta\}^2]^{1/2} \right. \\
 & \left. + 3\lambda_\beta\sqrt{s}/\alpha_0\beta, 8a\sqrt{s}\gamma\sqrt{\log\{\max(p, mn)\}/(mn)}, (4\alpha_0\beta)^{-1/2}\{2\sigma_{1F}^2\log\{\max(p, mn)\}/(mn)\}^{1/4} \right)
 \end{aligned}$$

Proof: First note that

$$\begin{aligned}
 & \left\| \frac{\partial \mathcal{L}(\boldsymbol{\Theta}_0, \widehat{\boldsymbol{\beta}})}{\partial \boldsymbol{\Theta}} - \frac{\partial \mathcal{L}(\boldsymbol{\Theta}_0, \boldsymbol{\beta}_0)}{\partial \boldsymbol{\Theta}} \right\|_{op} \\
 = & \left\| \frac{\partial^2 \mathcal{L}(\boldsymbol{\Theta}_0, \boldsymbol{\beta}^*)}{\partial \boldsymbol{\Theta} \partial \boldsymbol{\beta}^T} (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) \right\|_{op} \\
 = & \left\| (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m R_{ij} (H_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}^*) - E\{H_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}^*) | Y_{ij}\}) \right. \\
 & \left. + E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}^*)^{\otimes 2} | Y_{ij}\} - [E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}^*) | Y_{ij}\}]^{\otimes 2}) \mathbf{e}_i \mathbf{e}_j^T \{\mathbf{X}_{ij}^T (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)\} \right\|_{op} \\
 \leq & (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m \left\| R_{ij} (H_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}^*) - E\{H_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}^*) | Y_{ij}\}) \right. \\
 & \left. + E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}^*)^{\otimes 2} | Y_{ij}\} - [E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \boldsymbol{\Theta}_0, \boldsymbol{\beta}^*) | Y_{ij}\}]^{\otimes 2}) \mathbf{e}_i \mathbf{e}_j^T \right\|_{op} \max_{ij} |\mathbf{X}_{ij}^T (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)| \\
 \leq & (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m \left\| (2d_{H_2} + 2d_{S_2}^2) \mathbf{e}_i \mathbf{e}_j^T \right\|_{op} \max_{ij} |\mathbf{X}_{ij}^T (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)| \\
 = & (2d_{H_2} + 2d_{S_2}^2) \max \left( [320(d_{H_2} + d_{S_2}^2)a^2c_0^2/(\alpha_0\beta\sqrt{mn}) + \{3\lambda_\beta\sqrt{s}/\alpha_0\beta\}^2]^{1/2} \right. \\
 & \left. + 3\lambda_\beta\sqrt{s}/\alpha_0\beta, 8a\sqrt{s}\gamma\sqrt{\log\{\max(p, mn)\}/(mn)}, (4\alpha_0\beta)^{-1/2}\{2\sigma_{1F}^2\log\{\max(p, mn)\}/(mn)\}^{1/4} \right),
 \end{aligned}$$

where  $\boldsymbol{\beta}^*$  is on the line in between  $\widehat{\boldsymbol{\beta}}$  and  $\boldsymbol{\beta}_0$ . The last inequality holds by Condition (C5) and Theorem 1. This proves the result.

**Lemma A.15** *The  $\delta$ -covering number of  $\bar{\mathbb{B}}_{\boldsymbol{\Theta}}(\delta)$  satisfies*

$$\log N_{\boldsymbol{\Theta}}(\delta) \leq \frac{144\rho(\nu, D)^2}{\delta^2} \max(n, m).$$

Proof: First using the same arguments as those lead to (39) in Negahban and Wainwright (2012), we have

$$\sqrt{\log N_{\boldsymbol{\Theta}}(\delta)} \leq \frac{3\rho(\nu, D)}{\delta} E(\|\mathbf{G}\|_2),$$

where  $\mathbf{G}$  is a random matrix containing independent identically distribution standard normal entries. Further,

$$E(\|\mathbf{G}\|_2) \leq 4 \max(\sqrt{n}, \sqrt{m})$$

by the results in Section 3.1 in Bandeira et al. (2016). We have

$$\log N_{\boldsymbol{\Theta}}(\delta) \leq \frac{144\rho(\nu, D)^2}{\delta^2} \max(n, m).$$

This proves the result.



**Lemma A.16** *Assume Conditions (C1), (C5), (C6). Let  $\delta = D/\xi$ . Let  $\Theta^l \in \bar{\mathbb{B}}_{\Theta}(D)$ , and  $\mathbf{b}^l, l = 1, \dots, N_{\Theta}(\delta)$  be  $\delta$ -covering of  $\bar{\mathbb{B}}_{\Theta}(D)$  in Frobenius norm. Assume*

$$\frac{mn\delta^4}{128(d_{H_2} + d_{S_2}^2)a^2D^2} \geq 2\frac{144\rho(\nu, D)^2}{\delta^2} \max(n, m).$$

Recall that

$$\begin{aligned} & \text{vec}(\mathbf{A})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{B}) \\ = & (mn)^{-1} \left( \sum_{i=1}^n \sum_{j=1}^m \text{vec}(\mathbf{A})^T R_{ij} [H_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T \right. \\ & - E\{H_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \\ & \left. + E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}, \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} - E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) | Y_{ij}\}^{\otimes 2}] \text{vec}(\mathbf{B}) \right). \end{aligned}$$

Then there are positive constants  $c_{d2}, c_{d3}$  such that

$$\begin{aligned} & \Pr \left( \sup_{k=1, \dots, N(\delta)} \left| \text{vec}(\Theta^k)^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\Theta^k) \right. \right. \\ & \quad \left. \left. - \text{vec}(\Theta^k)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{z}_{ij} \mathbf{z}_{ij}^T | Y_{ij}\} \right. \right. \\ & \quad \left. \left. - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{z}_{ij} | Y_{ij}\}^{\otimes 2}] \text{vec}(\Theta^k) \right| \geq \delta^2 + \frac{32(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}} \right) \\ & \leq 4 \exp(-c_{d2} D^2 mn), \end{aligned}$$

$$\begin{aligned} & \Pr \left( \sup_{k=1, \dots, N(\delta)} \left| \text{vec}(\Theta^k)^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\Theta - \Theta^k) \right. \right. \\ & \quad \left. \left. - \text{vec}(\Theta^k)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{z}_{ij} \mathbf{z}_{ij}^T | Y_{ij}\} \right. \right. \\ & \quad \left. \left. - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \mathbf{z}_{ij} | Y_{ij}\}^{\otimes 2}] \text{vec}(\Theta - \Theta^k) \right| \geq \delta^2 + \frac{64(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}} \right) \\ & \leq 4 \exp(-c_{d3} D^2 mn). \end{aligned}$$

Proof: First we have

$$\begin{aligned} & \left| \text{vec}(\Theta)^{k^T} R_{ij} [H_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T - E\{H_2(Y_{ij}, \mathbf{X} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \right. \\ & \quad \left. + E\{S_2^2(Y_{ij}, \mathbf{X} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} - E\{S_2(Y_{ij}, \mathbf{X} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) | Y_{ij}\}^{\otimes 2}] \text{vec}(\Theta)^k \right| \\ & \leq \{2d_{H_2} \|\Theta^k\|_F \|\text{vec}(\mathbf{z}_{ij})\|_2 \|\Theta^k\|_{\max} \|\text{vec}(\mathbf{z}_{ij})\|_1 \\ & \quad + 2d_{S_2}^2 \|\Theta^k\|_F \|\text{vec}(\mathbf{z}_{ij})\|_2 \|\Theta^k\|_{\max} \|\text{vec}(\mathbf{z}_{ij})\|_1\} \\ & = 2(d_{H_2} + d_{S_2}^2)aD \end{aligned}$$

and

$$\left| \text{vec}(\Theta)^{k^T} R_{ij} [H_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T - E\{H_2(Y_{ij}, \mathbf{X} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \right.$$

$$\begin{aligned}
 & + E\{S_2^2(Y_{ij}, \mathbf{X}|\tilde{\Theta}, \tilde{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^\top|Y_{ij}\} - E\{S_2(Y_{ij}, \mathbf{X}|\tilde{\Theta}, \tilde{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2}\text{vec}(\Theta^k) \\
 \leq & \{2d_{H_2}\|\Theta^k\|_{\max}\|\text{vec}(\mathbf{z}_{ij})\|_1\|\Theta^k\|_{\max}\|\text{vec}(\mathbf{z}_{ij})\|_1 \\
 & + 2d_{S_2}^2\|\Theta^k\|_{\max}\|\text{vec}(\mathbf{z}_{ij})\|_1\|\Theta^k\|_{\max}\|\text{vec}(\mathbf{z}_{ij})\|_1\} \\
 = & 2(d_{H_2} + d_{S_2}^2)a^2.
 \end{aligned}$$

by Condition (C6). We now use Lemma A.2, where we treat each summand in  $\text{vec}(\Theta^k)^\top \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\text{vec}(\Theta^k)$  as  $Y_i$  in Lemma A.2, and set  $(u_i, v_i) = \{-2(d_{H_2} + d_{S_2}^2)aD, 2(d_{H_2} + d_{S_2}^2)aD\}$  and  $(u_{1i}, v_{1i}) = \{-2(d_{H_2} + d_{S_2}^2)a^2, 2(d_{H_2} + d_{S_2}^2)a^2\}$ . Consider  $\mathcal{T}$  that contains only one vector  $\mathbf{t}$ , where  $\mathbf{t}$  is the  $mn$  dimensional vector with element  $(mn)^{-1}$ . Then  $Z$  in Lemma A.2 is  $\text{vec}(\Theta^k)^\top \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\text{vec}(\Theta^k)$  and the  $\sigma$  in Lemma A.2 (7) is  $\sigma = [\sum_{i=1}^n \sum_{j=1}^m (mn)^{-2} \{4(d_{H_2} + d_{S_2}^2)aD\}^2]^{1/2} = 4(d_{H_2} + d_{S_2}^2)aD/\sqrt{mn} = O(1)$  by Condition (C6). Further, the  $\sigma$  in Lemma A.2 (8) is  $\sigma_1 = [\sum_{i=1}^n \sum_{j=1}^m (mn)^{-2} \{4(d_{H_2} + d_{S_2}^2)a^2\}^2]^{1/2} = 4(d_{H_2} + d_{S_2}^2)a^2/\sqrt{mn}$ . Hence let  $m_F$  be the median of  $\text{vec}(\Theta^k)^\top \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\text{vec}(\Theta^k)$ , we have

$$\begin{aligned}
 & \Pr\left(|\text{vec}(\Theta^k)^\top \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\text{vec}(\Theta^k) - \text{vec}(\Theta^k)^\top E\{\mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\}\text{vec}(\Theta^k)| \geq \delta^2 + \frac{32(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}}\right) \\
 \leq & \Pr\left(|\text{vec}(\Theta^k)^\top \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\text{vec}(\Theta^k) - m_F| + |\text{vec}(\Theta^k)^\top E\{\mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\}\text{vec}(\Theta^k) - m_F| \geq \delta^2\right. \\
 & \left. + \frac{32(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}}\right) \\
 \leq & \Pr\left(|\text{vec}(\Theta^k)^\top \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\text{vec}(\Theta^k) - m_F| \geq \delta^2 + \frac{32(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}} - \frac{16\sqrt{\pi}(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}}\right) \\
 \leq & \Pr\left(|\text{vec}(\Theta^k)^\top \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\text{vec}(\Theta^k) - m_F| \geq \delta^2\right) \\
 \leq & 4 \exp\left\{\frac{-mn\delta^4}{64(d_{H_2} + d_{S_2}^2)^2 a^2 D^2}\right\}.
 \end{aligned}$$

Now for  $\Theta^k \in \bar{\mathbb{B}}(D)$ ,  $k = 1, \dots, N_\Theta(\delta)$ , we have

$$\begin{aligned}
 & \Pr\left\{\sup_{k=1, \dots, N(\delta)} |\text{vec}(\Theta^k)^\top \mathbf{F}_\Theta(\Theta^k, \mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\text{vec}(\Theta^k) \right. \\
 & \quad - (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m \text{vec}(\Theta^k)^\top E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^\top|Y_{ij}\} \\
 & \quad \left. - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\tilde{\Theta}, \tilde{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2}\text{vec}(\Theta^k)| \geq \delta^2 + \frac{32(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}}\right\} \\
 \leq & 4 \exp\left\{\frac{-mn\delta^4}{64(d_{H_2} + d_{S_2}^2)^2 a^2 D^2} + \log\{N_\Theta(\delta)\}\right\} \\
 \leq & 4 \exp\left\{\frac{-mn\delta^4}{64(d_{H_2} + d_{S_2}^2)^2 a^2 D^2} + \frac{144\rho(\nu, D)^2}{\delta^2} \max(n, m)\right\}.
 \end{aligned}$$

Now recall  $\delta = D/\xi$ , and

$$\frac{mn\delta^4}{64(d_{H_2} + d_{S_2}^2)^2 a^2 D^2} \geq 2 \frac{144\rho(\nu, D)^2}{\delta^2} \max(n, m).$$

Thus we have

$$\Pr\left\{\sup_{k=1, \dots, N_\Theta(\delta)} |\text{vec}(\Theta^k)^\top \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\text{vec}(\Theta^k) \right.$$

$$\begin{aligned}
 & -\text{vec}(\Theta^k)^\top (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E\{R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^\top | Y_{ij}\} \\
 & - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) | Y_{ij}\}^{\otimes 2}\} \text{vec}(\Theta)^k \geq \delta^2 + \frac{32(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}} \Big\} \\
 & \leq 4 \exp(-c_{d2} D^2 mn),
 \end{aligned}$$

where  $c_{d2} = 1/\{128a^2\xi^4(d_{H_2} + d_{S_2}^2)^2\}$ . Using the similar argument we can show that there is a  $c_{d3} > 0$  such that

$$\begin{aligned}
 & \Pr \left\{ \sup_{k=1, \dots, N_{\Theta}(\delta)} |\text{vec}(\Theta^k)^\top \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\Theta - \Theta^k) \right. \\
 & - (\text{vec} \Theta^k)^\top (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E\{R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^\top | Y_{ij}\} \\
 & \left. - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) | Y_{ij}\}^{\otimes 2}\} \text{vec}(\Theta - \Theta^k) \geq \delta^2 + \frac{64(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}} \right\} \\
 & \leq 4 \exp(-c_{d3} D^2 mn).
 \end{aligned}$$

with  $c_{d3} = 1/\{256a^2\xi^4(d_{H_2} + d_{S_2}^2)^2\}$ . This proves the result.

Define

$$\begin{aligned}
 Q_{ij} & \equiv |H_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) - E\{H_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) | Y_{ij}\} \\
 & \quad + E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) | Y_{ij}\} - E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \tilde{\Theta}, \tilde{\beta}) | Y_{ij}\}^2|,
 \end{aligned} \tag{17}$$

and by Condition (C6), we have  $|Q_{ij}| \leq 2(d_{H_2} + d_{S_2}^2)$ .

**Lemma A.17** *Define*

$$\mathcal{D}(D) = \{\Delta_{\Theta} \in \mathbb{R}^{n \times m} \mid \|\Delta_{\Theta}\|_F \leq \delta, \|\Delta_{\Theta}\|_* \leq 2\rho(\nu, D), \text{ and } \|\Delta_{\Theta}\|_{\max} \leq 2a\}.$$

The  $\delta$ -covering number of  $\mathcal{D}(D)$  satisfies

$$\log N_{\mathcal{D}}(\delta) \leq \frac{576\rho(\nu, D)^2}{\delta^2} \max(n, m).$$

Proof: Because  $\mathcal{D}(D) \subset \tilde{\mathcal{D}}(D) \equiv \{\Delta_{\Theta} \in \mathbb{R}^{n \times m} \mid \|\Delta_{\Theta}\|_* \leq 2\rho(\nu, D)\}$ . Then  $\delta$ -covering number of  $\tilde{\mathcal{D}}(D)$  defined by  $N_{\tilde{\mathcal{D}}}(\delta)$  satisfies  $N_{\mathcal{D}}(\delta) \leq N_{\tilde{\mathcal{D}}}(\delta)$ . Using the same arguments as those lead to (39) in Negahban and Wainwright (2012), we have

$$\sqrt{\log N_{\tilde{\mathcal{D}}}(\delta)} \leq \frac{6\rho(\nu, D)}{\delta} E(\|\mathbf{G}\|_2),$$

where  $\mathbf{G}$  is a random matrix containing independent identically distribution standard normal entries. Further,

$$E(\|\mathbf{G}\|_2) \leq 4 \max(\sqrt{n}, \sqrt{m})$$

by the results in Section 3.1 in Bandeira et al. (2016). We have

$$\log N_{\mathcal{D}}(\delta) \leq \log N_{\tilde{\mathcal{D}}}(\delta) \leq \frac{576\rho(\nu, D)^2}{\delta^2} \max(n, m).$$

This proves the result.

**Lemma A.18** *Assume Conditions (C4), (C5) and (C6) hold. Define*

$$\mathcal{D}(D) = \{\Delta_{\Theta} \in \mathbb{R}^{n \times m} \mid \|\Delta_{\Theta}\|_F \leq \delta, \|\Delta_{\Theta}\|_* \leq 2\rho(\nu, D), \text{ and } \|\Delta_{\Theta}\|_{\max} \leq 2a\},$$

and let  $\delta = D/\xi$ . Assume  $\nu$  satisfies

$$\frac{mn\delta^4}{256(d_{H_2} + d_{S_2}^2)^2 a^2 \delta^2} \geq 2 \frac{576\rho(\nu, D)^2}{\delta^2} \max(n, m).$$

Further, let  $\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} \mid \tilde{\Theta}, \tilde{\beta})$  be as defined in Lemma A.16. Then

$$\begin{aligned} & \Pr \left\{ \sup_{\Delta_{\Theta} \in \mathcal{D}(D)} |\text{vec}(\Delta_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} \mid \tilde{\Theta}, \tilde{\beta}) \text{vec}(\Delta_{\Theta}) - \text{vec}(\Delta_{\Theta})^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} \mid \tilde{\Theta}, \tilde{\beta})\} \text{vec}(\Delta_{\Theta})| > D^2/\xi^2 \right. \\ & \quad \left. + \frac{114(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}} \right\} \\ & \leq 4 \exp \left( \frac{-mnD^2}{512(d_{H_2} + d_{S_2}^2)^2 a^2 \xi^2} \right). \end{aligned}$$

Proof: First note that for  $\Delta_{\Theta} \in \mathcal{D}(D)$  we have

$$\begin{aligned} & |\text{vec}(\Delta_{\Theta})^T R_{ij} H_2(Y_{ij}, \mathbf{X}_{ij} \mid \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T \text{vec}(\Delta_{\Theta})| \\ & \leq \{ \|\Delta_{\Theta}\|_{\max} \|\text{vec}(\mathbf{z}_{ij})\|_1 \}^2 H_2(Y_{ij}, \mathbf{X}_{ij}, \tilde{\Theta}, \tilde{\beta}) \leq 4a^2 d_{H_2}. \end{aligned}$$

Similarly,

$$\begin{aligned} |\text{vec}(\Delta_{\Theta})^T R_{ij} E\{H_2(Y_{ij}, \mathbf{X}_{ij} \mid \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \text{vec}(\Delta_{\Theta})| & \leq 4a^2 d_{H_2}, \\ |\text{vec}(\Delta_{\Theta})^T R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} \mid \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \text{vec}(\Delta_{\Theta})| & \leq 4a^2 d_{S_2}^2, \\ |\text{vec}(\Delta_{\Theta})^T R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} \mid \tilde{\Theta}, \tilde{\beta}) \mathbf{z}_{ij} | Y_{ij}\}^{\otimes 2} \text{vec}(\Delta_{\Theta})| & \leq 4a^2 d_{S_2}^2. \end{aligned}$$

Also we have

$$\begin{aligned} & |\text{vec}(\Delta_{\Theta})^T R_{ij} H_2(Y_{ij}, \mathbf{X}_{ij} \mid \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T \text{vec}(\Delta_{\Theta})| \\ & \leq \{ \|\Delta_{\Theta}\|_{\max} \|\text{vec}(\mathbf{z}_{ij})\|_1 \} \{ \|\Delta_{\Theta}\|_F \|\mathbf{z}_{ij}\|_F \} H_2(Y_{ij}, \mathbf{X}_{ij}, \Theta, \beta^*) \leq 2a\delta d_{H_2}. \end{aligned}$$

Similarly,

$$\begin{aligned} |\text{vec}(\Delta_{\Theta})^T R_{ij} E\{H_2(Y_{ij}, \mathbf{X}_{ij} \mid \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \text{vec}(\Delta_{\Theta})| & \leq 2a\delta d_{H_2}, \\ |\text{vec}(\Delta_{\Theta})^T R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} \mid \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \text{vec}(\Delta_{\Theta})| & \leq 2a\delta d_{S_2}^2, \\ |\text{vec}(\Delta_{\Theta})^T R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} \mid \tilde{\Theta}, \tilde{\beta}) \mathbf{z}_{ij} | Y_{ij}\}^{\otimes 2} \text{vec}(\Delta_{\Theta})| & \leq 2a\delta d_{S_2}^2. \end{aligned}$$

Hence each summand in  $\text{vec}(\Delta_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} \mid \tilde{\Theta}, \tilde{\beta}) \text{vec}(\Delta_{\Theta})$  is in the range of  $[-8a^2(d_{H_2} + d_{S_2}^2), 8a^2(d_{H_2} + d_{S_2}^2)]$ , and also in the range of  $[-4a\delta(d_{H_2} + d_{S_2}^2), 4a\delta(d_{H_2} + d_{S_2}^2)]$ . Define  $\sigma_d \equiv 8a\delta(d_{H_2} + d_{S_2}^2)$  be the  $\sigma$  in Lemma A.2 (7) and  $\sigma_{1d} \equiv 16a^2(d_{H_2} + d_{S_2}^2)$  be the  $\sigma$  in Lemma A.2 (8). By Lemma A.2, let  $m_F$  be the median of  $\text{vec}(\Delta_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} \mid \tilde{\Theta}, \tilde{\beta}) \text{vec}(\Delta_{\Theta})$ , we have

$$\begin{aligned} & \Pr \left( \left| \text{vec}(\Delta_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} \mid \tilde{\Theta}, \tilde{\beta}) \text{vec}(\Delta_{\Theta}) - \text{vec}(\Delta_{\Theta})^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} \mid \tilde{\Theta}, \tilde{\beta})\} \text{vec}(\Delta_{\Theta}) \right| \geq \right. \\ & \quad \left. \delta^2 + \frac{114(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}} \right) \end{aligned}$$

$$\begin{aligned}
 &\leq \Pr \left( |\text{vec}(\mathbf{\Delta}_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{\Delta}_{\Theta}) - m_F| + |\text{vec}(\mathbf{\Delta}_{\Theta})^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta})\} \text{vec}(\mathbf{\Delta}_{\Theta}) - m_F| \geq \right. \\
 &\quad \left. \delta^2 + \frac{114(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}} \right) \\
 &\leq \Pr \left( |\text{vec}(\mathbf{\Delta}_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{\Delta}_{\Theta}) - m_F| \geq \delta^2 \right. \\
 &\quad \left. + \frac{114(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}} - \frac{64\sqrt{\pi}(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}} \right) \\
 &\leq \Pr \left( \sqrt{mn} |\text{vec}(\mathbf{\Delta}_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{\Delta}_{\Theta}) - m_F| \geq \sqrt{mn}\delta^2 \right) \\
 &\leq 4 \exp \left\{ \frac{-mn\delta^4}{256(d_{H_2} + d_{S_2}^2)^2 a^2 \delta^2} \right\}.
 \end{aligned}$$

Therefore combine Lemma A.17,

$$\begin{aligned}
 &\Pr \left\{ \sup_{\mathbf{\Delta}_{\Theta} \in \mathcal{D}(D)} |\text{vec}(\mathbf{\Delta}_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{\Delta}_{\Theta}) - \text{vec}(\mathbf{\Delta}_{\Theta})^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta})\} \text{vec}(\mathbf{\Delta}_{\Theta})| > \delta^2 \right. \\
 &\quad \left. + \frac{114(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}} \right\} \\
 &\leq 4 \exp \left( \frac{-mn\delta^4}{256(d_{H_2} + d_{S_2}^2)^2 a^2 \delta^2} + \log\{N_{\mathcal{D}}(\delta)\} \right) \\
 &\leq 4 \exp \left( \frac{-mn\delta^4}{256(d_{H_2} + d_{S_2}^2)^2 a^2 \delta^2} + \frac{576\rho(\nu, D)^2}{\delta^2} \max(n, m) \right).
 \end{aligned}$$

Now because  $\delta = D/\xi$  and

$$\frac{mn\delta^4}{256(d_{H_2} + d_{S_2}^2)^2 a^2 \delta^2} \geq 2 \frac{576\rho(\nu, D)^2}{\delta^2} \max(n, m).$$

Therefore,

$$\begin{aligned}
 &\Pr \left\{ \sup_{\mathbf{\Delta}_{\Theta} \in \mathcal{D}(D)} |\text{vec}(\mathbf{\Delta}_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta}) \text{vec}(\mathbf{\Delta}_{\Theta}) - \text{vec}(\mathbf{\Delta}_{\Theta})^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta})\} \text{vec}(\mathbf{\Delta}_{\Theta})| > D^2/\xi^2 \right. \\
 &\quad \left. + \frac{114(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}} \right\} \\
 &\leq 4 \exp \left( \frac{-mnD^2}{512(d_{H_2} + d_{S_2}^2)^2 a^2 \xi^2} \right).
 \end{aligned}$$

This proves the result.

**Lemma A.19** *Assume Conditions (C1)–(C5) hold. Let  $\delta = D/\xi$ . Then we have*

$$\begin{aligned}
 &\Pr \left\{ \left| \text{vec}(\Theta)^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \Theta^*, \hat{\beta}) \text{vec}(\Theta) - \text{vec}(\Theta)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \Theta^*, \hat{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \right. \right. \\
 &\quad \left. \left. - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \Theta^*, \hat{\beta}) \text{vec}(\mathbf{z}_{ij}) | Y_{ij}\}^{\otimes 2}] \text{vec}(\Theta) \right| \right. \\
 &\geq \alpha_{\min} \left( (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \Theta^*, \hat{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \right.
 \end{aligned}$$

$$\begin{aligned}
 & -R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2}) / 2\|\Theta\|_F^2 + \frac{274(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}}, \Theta \in \mathcal{C}_\Theta(\nu) \} \\
 = & \exp\{-Cd\log(d)\}
 \end{aligned}$$

for some positive constant  $C$ , where the expectations are taken over  $\mathbf{X}, \mathbf{Y}$ .

Proof: For any  $\Theta \in \bar{\mathbb{B}}_\Theta(D)$ , for any  $k \in \{1, \dots, N(\delta)\}$ , let  $\widetilde{\Delta}_\Theta = \Theta - \Theta^k$ , then  $\widetilde{\Delta}_\Theta \in \mathcal{D}(D)$ . In addition,

$$\begin{aligned}
 & \text{vec}(\Theta)^T \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta}) \text{vec}(\Theta) - \text{vec}(\Theta)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^T|Y_{ij}\} \\
 & - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2}] \text{vec}(\Theta) \\
 = & \mathbf{F}_\beta(\Theta^k + \widetilde{\Delta}_\Theta, \Theta^k + \widetilde{\Delta}_\Theta, \mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta}) \\
 & - \text{vec}(\Theta^k + \widetilde{\Delta}_\Theta)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^T|Y_{ij}\} \\
 & - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2}] \text{vec}(\Theta^k + \widetilde{\Delta}_\Theta) \\
 = & \text{vec}(\Theta^k)^T \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta}) \text{vec}(\Theta^k) - \text{vec}(\Theta^k)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^T|Y_{ij}\} \\
 & - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2}] \text{vec}(\Theta^k) \\
 & + 2\text{vec}(\Theta^k)^T \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta}) \widetilde{\Delta}_\Theta - 2\text{vec}(\Theta^k)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^T|Y_{ij}\} \\
 & - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2}] \text{vec}(\widetilde{\Delta}_\Theta) \\
 & + \text{vec}(\widetilde{\Delta}_\Theta)^T \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta}) \widetilde{\Delta}_\Theta - \text{vec}(\widetilde{\Delta}_\Theta)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^T|Y_{ij}\} \\
 & - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2}] \text{vec}(\widetilde{\Delta}_\Theta).
 \end{aligned}$$

Hence by Lemmas A.16 and A.18, while replacing  $\widetilde{\Theta}$  by  $\Theta^*$  and  $\widetilde{\beta}$  by  $\widehat{\beta}$ , there are constants  $b_{d3}$  such that

$$\sup_{\Theta \in \bar{\mathbb{B}}_\Theta(D)} |\text{vec}(\Theta)^T \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta}) \text{vec}(\Theta)| \quad (18)$$

$$\begin{aligned}
 & - \text{vec}(\Theta)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^T|Y_{ij}\} \\
 & - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2}] \text{vec}(\Theta)| \\
 \leq & \sup_{\Theta \in \bar{\mathbb{B}}_\Theta(D)} |\text{vec}(\Theta)^T \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta}) \text{vec}(\Theta)| \quad (19)
 \end{aligned}$$

$$\begin{aligned}
 & - \text{vec}(\Theta)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^T|Y_{ij}\} \\
 & - R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2}] \text{vec}(\Theta)| \\
 \leq & \sup_{k=1, \dots, N_\Theta(\delta)} |\text{vec}(\Theta^k)^T \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta}) \text{vec}(\Theta^k)| \quad (20)
 \end{aligned}$$

$$- \text{vec}(\Theta^k)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^T|Y_{ij}\}$$

$$\begin{aligned}
 & -R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2}\text{vec}(\Theta^k)| \\
 & +2 \sup_{k=1, \dots, N_{\Theta}(\delta)} |\text{vec}(\Theta^k)^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta}) \widetilde{\Delta}_{\Theta}| \tag{21}
 \end{aligned}$$

$$\begin{aligned}
 & -\text{vec}(\Theta^k)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^T|Y_{ij}\} \\
 & -R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2}]\text{vec}(\widetilde{\Delta}_{\Theta})| \\
 & + \sup_{\widetilde{\Delta}_{\Theta} \in \mathcal{D}(D)} |\text{vec}(\widetilde{\Delta}_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta})\text{vec}(\widetilde{\Delta}_{\Theta})| \tag{22}
 \end{aligned}$$

$$\begin{aligned}
 & -\text{vec}(\widetilde{\Delta}_{\Theta})^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^T|Y_{ij}\} \\
 & -R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2}]\text{vec}(\widetilde{\Delta}_{\Theta})| \\
 & \leq 4D^2/\xi^2 + \frac{274(d_{H_2} + d_{S_2}^2)a^2}{\sqrt{mn}} \tag{23}
 \end{aligned}$$

with probability at least  $1 - 12\{\exp(-b_{d_3}D^2mn)\}$ . The last inequality holds by Lemma A.16 and A.18. Denote

$$\begin{aligned}
 \alpha_0 \equiv & \alpha_{\min} \left( (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^T|Y_{ij}\} \right. \\
 & \left. -R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2} \right),
 \end{aligned}$$

We choose  $\xi > \sqrt{8/\alpha_0}$  and define  $\alpha = \xi^2\alpha_0/8$ , then  $\alpha > 1$ .

Now note for  $\Theta \in \mathcal{C}_{\Theta}(\nu)$  with  $\|\Theta\|_{\max} = b$  for a constant  $b$ , we have

$$\|\Theta\|_F^2 \geq b\nu\|\Theta\|_* \sqrt{\frac{\log(d)}{d}} \geq b\nu\|\Theta\|_F \sqrt{\frac{\log(d)}{d}},$$

which implies  $\|\Theta\|_F \geq b\nu\sqrt{\log(d)/d}$ . Define  $\mu(b) = b^2\nu^2\log(d)/d$ . Moreover, we define

$$\mathcal{S}_l(b) \equiv \{\Theta \in \mathcal{C}_{\Theta}(\nu) \mid \|\Theta\|_{\max} = b, \sqrt{\alpha^{l-1}\mu(b)} \leq \|\Theta\|_F \leq \sqrt{\alpha^l\mu(b)}, \|\Theta\|_* \leq \rho(\nu, \sqrt{\alpha^l\mu(b)})\},$$

then  $\mathcal{S}_l(b) \subset \mathbb{B}_{\Theta}(\sqrt{\alpha^l\mu(b)})$ . For  $\Theta \in \mathcal{S}_l$   $D_{d\tau} = 274(d_{H_2} + d_{S_2}^2)a^2$ , we have

$$\begin{aligned}
 (\alpha_0/2)\|\Theta\|_F^2 + \frac{D_{d\tau}}{\sqrt{mn}} & \geq (\alpha_0/2)\alpha^{l-1}\mu(b) + \frac{D_{d\tau}}{\sqrt{mn}} \\
 & = \frac{4\alpha^l\mu(b)}{\xi^2} + \frac{D_{d\tau}}{\sqrt{mn}}.
 \end{aligned}$$

Therefore,

$$\begin{aligned}
 & \Pr \left\{ |\text{vec}(\Theta)^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta})\text{vec}(\Theta) \right. \\
 & -\text{vec}(\Theta)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^T|Y_{ij}\} \\
 & \left. -R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2}]\text{vec}(\Theta) \right|
 \end{aligned}$$

$$\begin{aligned}
 &\geq (\alpha_0/2)\|\Theta\|_F^2 + \frac{D_{d7}}{\sqrt{mn}}, \Theta \in \mathcal{C}_\Theta(\nu) \Big\} \\
 \leq &\int_0^a \Pr \left\{ |\text{vec}(\Theta)^T \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta}) \text{vec}(\Theta) \right. \\
 &\quad - \text{vec}(\Theta)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \\
 &\quad \left. - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta}) \text{vec}(\mathbf{z}_{ij}) | Y_{ij}\}^{\otimes 2}] \text{vec}(\Theta) | \right. \\
 &\quad \left. \geq (\alpha_0/2)\|\Theta\|_F^2 + \frac{D_{d7}}{\sqrt{mn}}, \Theta \in \mathcal{C}_\Theta(\nu), \|\Theta\|_{\max} = b \right\} db \\
 \leq &\int_0^a \sum_{l=1}^{\infty} \Pr \left\{ |\text{vec}(\Theta)^T \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta}) \text{vec}(\Theta) \right. \\
 &\quad - \text{vec}(\Theta)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \\
 &\quad \left. - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta}) \text{vec}(\mathbf{z}_{ij}) | Y_{ij}\}^{\otimes 2}] \text{vec}(\Theta) | \right. \\
 &\quad \left. \geq (\alpha_0/2)\|\Theta\|_F^2 + \frac{D_{d7}}{\sqrt{mn}}, \Theta \in \mathcal{S}_l(b) \right\} db \\
 = &\sum_{l=1}^{\infty} \Pr \left\{ |\text{vec}(\Theta)^T \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta}) \text{vec}(\Theta) \right. \\
 &\quad - \text{vec}(\Theta)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \\
 &\quad \left. - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta}) \text{vec}(\mathbf{z}_{ij}) | Y_{ij}\}^{\otimes 2}] \text{vec}(\Theta) | \right. \\
 &\quad \left. \geq (\alpha_0/2)\|\Theta\|_F^2 + \frac{D_{d7}}{\sqrt{mn}}, \Theta \in \mathcal{S}_l(b^*) \right\} \\
 \leq &\sum_{l=1}^{\infty} \Pr \left\{ |\text{vec}(\Theta)^T \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\Theta^*, \widehat{\beta}) \text{vec}(\Theta) \right. \\
 &\quad - \text{vec}(\Theta)^T (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \\
 &\quad \left. - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \widehat{\beta}) \text{vec}(\mathbf{z}_{ij}) | Y_{ij}\}^{\otimes 2}] \text{vec}(\Theta) | \geq 4\alpha^l \mu(b^*)/\xi^2 + \frac{D_{d7}}{\sqrt{mn}}, \Theta \in \mathcal{S}_l(b^*) \right\} \\
 \leq &12 \sum_{l=1}^{\infty} \exp(-b_{d3} \mu(b^*) \alpha^l mn) \\
 \leq &12 \sum_{l=1}^{\infty} \exp\{-b_{d3} l \log(\alpha) \mu(b^*) mn\} \\
 \leq &12 \frac{\exp\{-b_{d3} \log(\alpha) mn \mu(b^*)\}}{1 - \exp\{-b_{d3} \log(\alpha) mn \mu(b^*)\}} \\
 \leq &\exp\{-Cd \log(d)\},
 \end{aligned}$$

for some positive constant  $C$ , where  $b^*$  is a point on the line connecting 0 and  $a$ . The third inequality holds by (18), the fourth inequality holds because  $\alpha^l \geq l \log(\alpha) > 0$  for  $\alpha > 1$ . This proves the result.



**Lemma A.20** Assume Conditions (C5) and (C6) hold, and select  $\lambda_{\Theta} \geq 2\|\partial\mathcal{L}(\Theta_0, \hat{\beta})/\partial\Theta\|_{op}$ . If

$$\text{vec}(\Delta_{\Theta})^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \hat{\beta})\} \text{vec}(\Delta_{\Theta}) \leq \sigma_{dF} \sqrt{2d\log(d)/(mn)},$$

then

$$\begin{aligned} \|\Delta_{\Theta}\|_F &\leq \alpha_{\min} \left( (mn)^{-1} \sum_{i=1}^n \sum_{j=1}^m E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \hat{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \right. \\ &\quad \left. - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \hat{\beta}) \text{vec}(\mathbf{z}_{ij}) | Y_{ij}\}^{\otimes 2}] \right)^{-1/2} \{2\sigma_{dF}^2 d\log(d)/(mn)\}^{1/4}. \end{aligned}$$

Otherwise

$$\|\Delta_{\Theta}\|_* \leq 8\sqrt{r}\|\Delta_{\Theta}\|_F$$

with probability at least  $1 - 2\exp\{-d\log(d)\}$ , where  $\sigma_{dF} = 32a^2(d_{H_2} + d_{S_2}^2)$  as defined in Theorem 2.

Proof: First as shown in Lemma A.18 that for  $\Delta_{\Theta} \in \mathcal{D}(D)$  we have

$$\begin{aligned} &|\text{vec}(\Delta_{\Theta})^T R_{ij} H_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \hat{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T \text{vec}(\Delta_{\Theta})| \\ &\leq \{\|\Delta_{\Theta}\|_{\max} \|\text{vec}(\mathbf{z}_{ij})\|_1\}^2 |H_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \hat{\beta})| \leq 4a^2 d_{H_2}. \end{aligned}$$

Similarly,

$$\begin{aligned} |\text{vec}(\Delta_{\Theta})^T R_{ij} E\{H_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \hat{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \text{vec}(\Delta_{\Theta})| &\leq 4a^2 d_{H_2}, \\ |\text{vec}(\Delta_{\Theta})^T R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \hat{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \text{vec}(\Delta_{\Theta})| &\leq 4a^2 d_{S_2}^2, \\ |\text{vec}(\Delta_{\Theta})^T R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \hat{\beta}) \mathbf{z}_{ij} | Y_{ij}\}^{\otimes 2} \text{vec}(\Delta_{\Theta})| &\leq 4a^2 d_{S_2}^2. \end{aligned}$$

Hence each summand in  $\text{vec}(\Delta_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \hat{\beta}) \text{vec}(\Delta_{\Theta}) - \text{vec}(\Delta_{\Theta})^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \hat{\beta})\} \text{vec}(\Delta_{\Theta})$  is in the range of  $[-16a^2(d_{H_2} + d_{S_2}^2), 16a^2(d_{H_2} + d_{S_2}^2)]$ , and in turn is sub-Gaussian with parameter  $\sigma_{dF}$  by Lemma A.1. Lemma A.1 further leads to

$$\begin{aligned} &\Pr \left\{ |\text{vec}(\Delta_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \hat{\beta}) \text{vec}(\Delta_{\Theta}) - \text{vec}(\Delta_{\Theta})^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \hat{\beta})\} \text{vec}(\Delta_{\Theta})| > t \right\} \\ &= \Pr \left\{ |mn \text{vec}(\Delta_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \hat{\beta}) \text{vec}(\Delta_{\Theta}) - mn \text{vec}(\Delta_{\Theta})^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \hat{\beta})\} \text{vec}(\Delta_{\Theta})| > mnt \right\} \\ &\leq 2 \exp \left\{ -\frac{(mn)^2 t^2}{2mn\sigma_{dF}^2} \right\} \\ &= 2 \exp \left( -\frac{mnt^2}{2\sigma_{dF}^2} \right) \end{aligned} \tag{24}$$

for any  $t > 0$ . By (24), let  $t = \sigma_{dF} \sqrt{2d\log(d)/(mn)}$ , we have

$$\text{vec}(\Delta_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \hat{\beta}) \text{vec}(\Delta_{\Theta}) \geq \text{vec}(\Delta_{\Theta})^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \hat{\beta})\} \text{vec}(\Delta_{\Theta}) - \sigma_{dF} \sqrt{2d\log(d)/(mn)} \tag{25}$$

with probability at least  $1 - 2\exp\{-d\log(d)\}$ .

Now we discuss two cases. Case I:

$$\text{vec}(\Delta_{\Theta})^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \hat{\beta})\} \text{vec}(\Delta_{\Theta}) \leq \sigma_{dF} \sqrt{2d\log(d)/(mn)}$$

This leads to

$$\alpha_{\min} \left( E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \hat{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \right.$$

$$-R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \hat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2})\|\Delta_{\Theta}\|_F^2 \leq \sigma_{dF}\sqrt{2d\log(d)/(mn)},$$

hence

$$\begin{aligned} \|\Delta_{\Theta}\|_F &\leq \alpha_{\min}(E[R_{ij}E\{S_2^2(Y_{ij}, \mathbf{X}_{ij}|\Theta^*, \hat{\beta})\text{vec}(\mathbf{z}_{ij})\text{vec}(\mathbf{z}_{ij})^T|Y_{ij}\} \\ &\quad -R_{ij}E\{S_2(Y_{ij}, \mathbf{X}_{ij}, \Theta^*, \hat{\beta})\text{vec}(\mathbf{z}_{ij})|Y_{ij}\}^{\otimes 2})]^{-1/2}\{2\sigma_{dF}^2d\log(d)/(mn)\})^{1/4}. \end{aligned}$$

Case II:

$$\text{vec}(\Delta_{\Theta})^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \hat{\beta})\}\text{vec}(\Delta_{\Theta}) > \sigma_{dF}\sqrt{2d\log(d)/(mn)}.$$

Under (25), this leads to

$$\text{vec}(\Delta_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \hat{\beta})\text{vec}(\Delta_{\Theta}) > 0. \quad (26)$$

By Lemma 1 and (21) in Negahban and Wainwright (2012), when  $\lambda_{\Theta} \geq 2\|\partial\mathcal{L}(\Theta_0, \hat{\beta})/\partial\beta\|_{op}$ , note that  $\Theta_0$  has rank  $r$ , which satisfies  $r \leq 2r$ , we get

$$\|\Delta_{\Theta}\|_* \leq 8\sqrt{r}\|\Delta_{\Theta}\|_F.$$

This proves the result.

## Appendix F. Proof of Theorem 2

Proof: Let  $\Delta_{\Theta} = \hat{\Theta} - \Theta_0$ , we consider the situation where

$$\text{vec}(\Delta_{\Theta})^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y}|\Theta^*, \hat{\beta})\}\text{vec}(\Delta_{\Theta}) > \sigma_{dF}\sqrt{2d\log(d)/(mn)}.$$

First we consider Case I:  $\Delta_{\Theta} \notin \mathcal{C}_{\Theta}(\nu)$ . Then, with probability at least  $1 - 2\exp\{-d\log(d)\}$ ,

$$\|\Delta_{\Theta}\|_F^2 \leq \|\Delta_{\Theta}\|_{\max}\|\Delta_{\Theta}\|_*\nu\sqrt{\frac{d\log(d)}{mn}} \leq 16a\sqrt{r}\|\Delta_{\Theta}\|_F\nu\sqrt{\frac{d\log(d)}{mn}},$$

which implies  $\|\Delta_{\Theta}\|_F \leq 16a\sqrt{r}\nu\sqrt{d\log(d)/(mn)}$ .

Case II:  $\Delta_{\Theta} \in \mathcal{C}_{\Theta}(\nu)$ . Because  $\hat{\Theta}, \hat{\beta}$  is the minimizer of  $\mathcal{L}(\Theta, \beta) + \lambda_{\Theta}\|\Theta\|_* + \lambda_{\beta}\|\beta\|_1$ , we have

$$\mathcal{L}(\hat{\Theta}, \hat{\beta}) - \mathcal{L}(\Theta_0, \hat{\beta}) = \left\langle \frac{\partial\mathcal{L}(\Theta_0, \hat{\beta})}{\partial\Theta}, \Delta_{\Theta} \right\rangle + \frac{1}{2}\text{vec}(\hat{\Theta} - \Theta_0)^T \frac{\partial\mathcal{L}(\Theta^*, \hat{\beta})}{\partial\text{vec}(\Theta)\partial\text{vec}(\Theta)^T} \text{vec}(\hat{\Theta} - \Theta_0)$$

and

$$\begin{aligned} \mathbf{F}_{\beta}(\Delta_{\Theta}, \Delta_{\Theta}, \mathbf{X}, \mathbf{Y}|\Theta^*, \hat{\beta})/2 &= \text{vec}(\hat{\Theta} - \Theta_0)^T \frac{\partial^2\mathcal{L}(\Theta^*, \hat{\beta})}{\partial\text{vec}(\Theta)\partial\text{vec}(\Theta)^T} \text{vec}(\hat{\Theta} - \Theta_0)/2 \\ &= \mathcal{L}(\hat{\Theta}, \hat{\beta}) - \mathcal{L}(\Theta_0, \hat{\beta}) - \left\langle \frac{\partial\mathcal{L}(\Theta_0, \hat{\beta})}{\partial\Theta}, \Delta_{\Theta} \right\rangle \\ &\leq \left\langle -\frac{\partial\mathcal{L}(\Theta_0, \hat{\beta})}{\partial\Theta}, \Delta_{\Theta} \right\rangle + \lambda_{\Theta}\|\Theta_0\|_* - \lambda_{\Theta}\|\hat{\Theta}\|_* \\ &\leq \left\| \frac{\partial\mathcal{L}(\Theta_0, \hat{\beta})}{\partial\Theta} \right\|_{op}\|\Delta_{\Theta}\|_* + \lambda_{\Theta}\|\Theta_0\|_* - \lambda_{\Theta}\|\hat{\Theta}\|_* \\ &\leq \lambda_{\Theta}/2\|\Delta_{\Theta}\|_* + \lambda_{\Theta}\|\Theta_0 - \hat{\Theta}\|_* \\ &\leq 12\lambda_{\Theta}\sqrt{r}\|\Delta_{\Theta}\|_F \end{aligned} \quad (27)$$

with probability at least  $1 - 2 \exp\{-d \log(d)\} - d^{-1}$ . The first inequality holds by the second order mean value theorem for  $\mathcal{L}(\widehat{\Theta}, \beta)$  on  $\Theta$ . The second inequality is a known property concerning norms and trace. The third inequality holds by the selection that

$$\begin{aligned} \lambda_{\Theta} \geq & 2c_d \sqrt{d \log(d)/(mn)} + 2(2d_{H_2} + 2d_{S_2}^2) \max \left( [10\sigma_{1F}/(\alpha_{0\beta} \sqrt{mn}) + \{3\lambda_{\beta} \sqrt{s}/\alpha_{0\beta}\}^2]^{1/2} \right. \\ & \left. + 3\lambda_{\beta} \sqrt{s}/\alpha_{0\beta}, 8a\sqrt{s}\gamma \sqrt{\log\{\max(p, mn)\}/(mn)}, (4\alpha_{0\beta})^{-1/2} \{2\sigma_{1F}^2 \log\{\max(p, mn)\}/(mn)\}^{1/4} \right), \end{aligned}$$

in the theorem statement which is greater than  $2\|\partial \mathcal{L}(\Theta_0, \widehat{\beta})/\partial \Theta\|_{op}$  with probability at least  $1 - d^{-1}$ . The last line holds by Lemma A.20.

Further, by Lemma A.19 we have

$$\begin{aligned} & \text{vec}(\Delta_{\Theta})^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \Theta^*, \widehat{\beta}) \text{vec}(\Delta_{\Theta})/2 \\ \geq & \alpha_{\min} (E[R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \Theta^*, \widehat{\beta}) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T | Y_{ij}\} \\ & - R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \Theta^*, \widehat{\beta}) \text{vec}(\mathbf{z}_{ij}) | Y_{ij}\}^{\otimes 2}]) / 4 \|\Delta_{\Theta}\|_F^2 - 137(d_{H_2} + d_{S_2}^2)a^2/\sqrt{mn} \\ = & \alpha_{0\Theta} \|\Delta_{\Theta}\|_F^2 - 137(d_{H_2} + d_{S_2}^2)a^2/\sqrt{mn} \end{aligned}$$

with probability at least  $1 - \exp\{-Cd \log(d)\}$ . Combine with (27), we have that with probability at least  $1 - \exp\{-Cd \log(d)\} - 2 \exp\{-d \log(d)\}$ ,

$$\alpha_{0\Theta} \|\Delta_{\Theta}\|_F^2 \leq 12\lambda_{\Theta} \sqrt{r} \|\Delta_{\Theta}\|_F + 137(d_{H_2} + d_{S_2}^2)a^2/(\sqrt{mn}).$$

Then

$$\{\|\Delta_{\Theta}\|_F - 6\lambda_{\Theta} \sqrt{r}/\alpha_{0\Theta}\}^2 \leq [137(d_{H_2} + d_{S_2}^2)a^2/(\alpha_{0\Theta} \sqrt{mn}) + 36\lambda_{\Theta}^2 r/\alpha_{0\Theta}^2]$$

Hence

$$\|\Delta_{\Theta}\|_F \leq [137(d_{H_2} + d_{S_2}^2)a^2/(\alpha_{0\Theta} \sqrt{mn}) + 36\lambda_{\Theta}^2 r/\alpha_{0\Theta}^2]^{1/2} + 6\lambda_{\Theta} \sqrt{r}/\alpha_{0\Theta}.$$

Combine with the order in Case I and Lemma A.20, we have

$$\begin{aligned} & \|\Delta_{\Theta}\|_F \\ \leq & \max \left( [137(d_{H_2} + d_{S_2}^2)a^2/(\alpha_{0\Theta} \sqrt{mn}) + 36\lambda_{\Theta}^2 r/\alpha_{0\Theta}^2]^{1/2} + 6\lambda_{\Theta} \sqrt{r}/\alpha_{0\Theta}, 16a\sqrt{r}\nu \sqrt{d \log(d)/mn}, \right. \\ & \left. (4\alpha_{0\Theta})^{-1/2} \{2\sigma_{dF}^2 d \log(d)/(mn)\}^{1/4} \right), \end{aligned}$$

with probability at least  $1 - \exp\{-Cd \log(d)\} - 2 \exp\{-d \log(d)\} - d^{-1}$ . This proves the result.

## Appendix G. Lemmas for Theorem 3

**Lemma A.21** *Assume Conditions (C5) and (C6) hold. Then for unit vectors  $\mathbf{v} \in \mathbb{R}^p$ ,  $\mathbf{u} \in \mathbb{R}^{mn}$  and  $\mathbf{w} \in \mathbb{R}^{mn+p}$  such that  $\|\mathbf{v}\|_2 = 1$ ,  $\|\mathbf{u}\|_2 = 1$  and  $\|\mathbf{w}\|_2 = 1$ , we have*

$$\mathbf{F}_{\beta}(\mathbf{v}, \mathbf{v}, \mathbf{X}, \mathbf{Y} | \Theta, \beta) \geq E\{\mathbf{F}_{\beta}(\mathbf{v}, \mathbf{v}, \mathbf{X}, \mathbf{Y} | \Theta, \beta)\} - 4c_0 \sqrt{(d_{H_2} + d_{S_2}^2) \log\{\max(p, mn)\}/(mn)}, \quad (28)$$

with probability at least  $1 - 2 \max(p, mn)^{-1}$ ,

$$\mathbf{u}^T \mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \Theta, \beta) \mathbf{u} \geq \mathbf{u}^T E\{\mathbf{F}_{\Theta}(\mathbf{X}, \mathbf{Y} | \Theta, \beta)\} \mathbf{u} - 4\sqrt{(d_{H_2} + d_{S_2}^2) \log(mn)/(mn)}, \quad (29)$$

with probability at least  $1 - 2(mn)^{-1}$ , and

$$\mathbf{w}^T \partial^2 \mathcal{L}(\Theta, \beta) / \partial \{\text{vec}(\Theta)^T, \beta^T\}^T \otimes^2 \mathbf{w}$$

$$\geq E\left(\mathbf{w}^T \partial^2 \mathcal{L}(\Theta, \beta) / \partial[\{\text{vec}(\Theta)^T, \beta^T\}^T]^{\otimes 2} \mathbf{w}\right) - 4(c_0 + 1) \sqrt{(d_{H_2} + d_{S_2}^2) \log(mn + p) / mn}$$

with probability at least  $1 - 2(p + mn)^{-1}$ .

Proof: First note that for unit vector  $\mathbf{v}$ ,  $\mathbf{v} \in \mathbb{R}^p$ , because  $\|\mathbf{v}\|_\infty \leq \|\mathbf{v}\|_2 = 1$ , we have

$$\begin{aligned} & |\mathbf{v}^T R_{ij} H_2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \mathbf{X}_{ij} \mathbf{X}_{ij}^T \mathbf{v}| \\ & \leq \{\|\mathbf{v}\|_\infty \|\mathbf{X}_{ij}\|_1\}^2 H_2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \leq c_0^2 d_{H_2}, \end{aligned}$$

Similarly

$$|\mathbf{v}^T R_{ij} E\{H_2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \mathbf{X}_{ij} \mathbf{X}_{ij}^T | Y_{ij}\} \mathbf{v}| \leq c_0^2 d_{H_2},$$

and

$$|\mathbf{v}^T R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \mathbf{X}_{ij} \mathbf{X}_{ij}^T | Y_{ij}\} \mathbf{v}| \leq c_0^2 d_{S_2}^2,$$

and

$$|\mathbf{v}^T R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \mathbf{X}_{ij} | Y_{ij}\}^{\otimes 2} \mathbf{v}| \leq c_0^2 d_{S_2}^2.$$

Hence we have each summand in  $\mathbf{F}_\beta(\mathbf{v}, \mathbf{v}, \mathbf{X}, \mathbf{Y} | \Theta, \beta) - E\{\mathbf{F}_\beta(\mathbf{v}, \mathbf{v}, \mathbf{X}, \mathbf{Y} | \Theta, \beta)\}$  is in the range of  $\{-4c_0^2(d_{H_2} + d_{S_2}^2), 4c_0^2(d_{H_2} + d_{S_2}^2)\}$ , and in turn is sub-Gaussian with parameter  $8c_0^2(d_{H_2} + d_{S_2}^2)$  by Lemma A.1. Here the expectation is taken over  $\mathbf{X}$  and  $\mathbf{Y}$ . Therefore, by Lemma A.1, for any given  $\mathbf{v}$

$$\begin{aligned} & \Pr\{|\mathbf{F}_\beta(\mathbf{v}, \mathbf{v}, \mathbf{X}, \mathbf{Y} | \Theta, \beta) - E\{\mathbf{F}_\beta(\mathbf{v}, \mathbf{v}, \mathbf{X}, \mathbf{Y} | \Theta, \beta)\}| > t\} \\ & = \Pr\{|mnF(\mathbf{v}, \mathbf{v}, \mathbf{X}, \mathbf{Y} | \Theta, \beta) - mnE\{\mathbf{F}_\beta(\mathbf{v}, \mathbf{v}, \mathbf{X}, \mathbf{Y} | \Theta, \beta)\}| > mnt\} \\ & \leq 2 \exp\left[-\frac{(mn)^2 t^2}{2mn\{8c_0^2(d_{H_2} + d_{S_2}^2)\}}\right] \\ & = 2 \exp\left\{-\frac{mnt^2}{16c_0^2(d_{H_2} + d_{S_2}^2)}\right\}. \end{aligned}$$

Now choose  $t = \sqrt{16c_0^2(d_{H_2} + d_{S_2}^2) \log\{\max(p, mn)\} / (mn)}$ , we obtain that for any given unit vector  $\mathbf{v}$ ,

$$\begin{aligned} & \Pr\left\{|\mathbf{F}_\beta(\mathbf{v}, \mathbf{v}, \mathbf{X}, \mathbf{Y} | \Theta, \beta) - E\{\mathbf{F}_\beta(\mathbf{v}, \mathbf{v}, \mathbf{X}, \mathbf{Y} | \Theta, \beta)\}| \right. \\ & \left. > \sqrt{16c_0^2(d_{H_2} + d_{S_2}^2) \log\{\max(p, mn)\} / (mn)}\right\} \leq 2 \max(p, mn)^{-1}. \end{aligned} \quad (30)$$

Therefore,

$$|\mathbf{F}_\beta(\mathbf{v}, \mathbf{v}, \mathbf{X}, \mathbf{Y} | \Theta, \beta) - E\{\mathbf{F}_\beta(\mathbf{v}, \mathbf{v}, \mathbf{X}, \mathbf{Y} | \Theta, \beta)\}| \leq 4 \sqrt{c_0^2(d_{H_2} + d_{S_2}^2) \log\{\max(p, mn)\} / (mn)},$$

with probability at least  $1 - 2 \max(p, mn)^{-1}$ . This proves (28).

Furthermore, for unit vector  $\mathbf{u} \in \mathbb{R}^{mn}$ , we have

$$\begin{aligned} & |\mathbf{u}^T R_{ij} H_2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^T \mathbf{u}| \\ & \leq \{\|\mathbf{u}\|_\infty \|\text{vec}(\mathbf{z}_{ij})\|_1\}^2 |H_2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta)| \leq d_{H_2}. \end{aligned}$$

Similarly,

$$\begin{aligned} |\mathbf{u}^\top R_{ij} E\{H_2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^\top | Y_{ij}\} \mathbf{u}| &\leq d_{H_2}, \\ |\mathbf{u}^\top R_{ij} E\{S_2^2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \text{vec}(\mathbf{z}_{ij}) \text{vec}(\mathbf{z}_{ij})^\top | Y_{ij}\} \mathbf{u}| &\leq d_{S_2}^2, \\ |\mathbf{u}^\top R_{ij} E\{S_2(Y_{ij}, \mathbf{X}_{ij} | \Theta, \beta) \text{vec}(\mathbf{z}_{ij}) | Y_{ij}\}^{\otimes 2} \mathbf{u}| &\leq d_{S_2}^2. \end{aligned}$$

Hence each summand in  $\mathbf{u}^\top \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y} | \Theta, \beta) \mathbf{u} - \mathbf{u}^\top E\{\mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y} | \Theta, \beta)\} \mathbf{u}$  is in the range of  $[-4(d_{H_2} + d_{S_2}^2), 4(d_{H_2} + d_{S_2}^2)]$ , and in turn is sub-Gaussian with parameter  $8(d_{H_2} + d_{S_2}^2)$  by Lemma A.1. Lemma A.1 further leads to

$$\begin{aligned} &\Pr\{|\mathbf{u}^\top \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y} | \Theta, \beta) \mathbf{u} - \mathbf{u}^\top E\{\mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y} | \Theta, \beta)\} \mathbf{u}| > t\} \\ &= \Pr\{|m n \mathbf{u}^\top \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y} | \Theta, \beta) \mathbf{u} - m n E\{\mathbf{u}^\top \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y} | \Theta, \beta)\} \mathbf{u}| > m n t\} \\ &\leq 2 \exp\left[-\frac{(m n)^2 t^2}{2 m n \{8(d_{H_2} + d_{S_2}^2)\}}\right] \\ &= 2 \exp\left\{-\frac{m n t^2}{16(d_{H_2} + d_{S_2}^2)}\right\} \end{aligned}$$

for any  $t > 0$ . Let  $t = \sqrt{16(d_{H_2} + d_{S_2}^2) \log(mn)/(mn)}$ , we have

$$\mathbf{u}^\top \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y} | \Theta, \beta) \mathbf{u} \geq \mathbf{u}^\top E\{\mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y} | \Theta, \beta)\} \mathbf{u} - 4\sqrt{(d_{H_2} + d_{S_2}^2) \log(mn)/(mn)} \quad (31)$$

with probability at least  $1 - 2(mn)^{-1}$ .

Moreover, using the same arguments as those lead to (30), (31), and the fact that  $\|\mathbf{z}_{ij}\|_1 = 1$ , we have for unit vector  $\mathbf{w} \in \mathbb{R}^{mn+p}$ ,

$$\begin{aligned} &\mathbf{w}^\top \partial^2 \mathcal{L}(\Theta, \beta) / \partial[\{\text{vec}(\Theta)^\top, \beta^\top\}^\top]^{\otimes 2} \mathbf{w} \\ &\geq E\left(\mathbf{w}^\top \partial^2 \mathcal{L}(\Theta, \beta) / \partial[\{\text{vec}(\Theta)^\top, \beta^\top\}^\top]^{\otimes 2} \mathbf{w}\right) - 4\sqrt{(c_0 + 1)^2 (d_{H_2} + d_{S_2}^2) \log(mn + p)/mn}, \end{aligned}$$

with probability at least  $1 - 2(p + mn)^{-1}$ .

**Lemma A.22** *Assume*

$$\beta^\top E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y} | \tilde{\Theta}, \tilde{\beta})\} \beta > 4c_0 \sqrt{(d_{H_2} + d_{S_2}^2) \log\{\max(p, mn)\}/(mn)} \|\beta\|_2^2,$$

for all  $\beta$  that satisfies  $\|\beta\|_\infty \leq 2a$  and  $\tilde{\beta}, \tilde{\Theta}$  that satisfy  $\|\tilde{\beta}\|_\infty \leq a$  and  $\|\tilde{\Theta}\|_{\max} \leq a$ . Recall that  $\sigma_\beta = C_1(2d_{H_2} + 2d_{S_2}^2)c_0^2$  for  $C_1 > 1$ ,  $Q^{t-1} = F(\Theta^{t-1}, \beta^{t-1})/\lambda_\beta$  and

$$g_\beta(\Theta^{t-1}, \beta^{t-1}, Q^{t-1}) = \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^{t-1}) / \partial \beta, \beta^{t-1} - \tilde{\beta}^{t-1} \rangle + \lambda_\beta \|\beta^{t-1} - \tilde{\beta}^{t-1}\|_1,$$

where

$$\tilde{\beta}^{t-1} \equiv \arg \min_\beta \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^{t-1}) / \partial \beta, \beta \rangle + \lambda_\beta \|\beta\|_1, \|\beta\|_1 \leq Q^{t-1}.$$

Select

$$0 < \eta < 1/\sigma_\beta.$$

Then we have

$$F(\Theta^{t-1}, \beta^t) \leq F(\Theta^{t-1}, \beta^{t-1}) - \frac{\eta g_\beta(\Theta^{t-1}, \beta^{t-1}, Q^{t-1})^2}{2(2Q^{t-1})^2},$$

with probability at least  $1 - 2 \max(p, mn)^{-1}$ .

Proof: Because

$$\boldsymbol{\beta}^\top E\{\mathbf{F}_\beta(\mathbf{X}, \mathbf{Y}|\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\beta}})\}\boldsymbol{\beta} > 4c_0\sqrt{(d_{H_2} + d_{S_2}^2)\log\{\max(p, mn)\}/(mn)}\|\boldsymbol{\beta}\|_2^2,$$

the second derivative of  $\mathcal{L}(\boldsymbol{\Theta}, \boldsymbol{\beta})$  with respect to  $\boldsymbol{\beta}$  is positive definite with probability at least  $1 - 2\max(p, mn)^{-1}$  by Lemma A.21, and hence we majorize  $F(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta})$  at  $\boldsymbol{\beta}^{t-1}$  as

$$\begin{aligned} F(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}) &\leq F(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1}) + \langle \boldsymbol{\beta} - \boldsymbol{\beta}^{t-1}, \partial\mathcal{L}(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1})/\partial\boldsymbol{\beta} \rangle \\ &\quad + \frac{\sigma_\beta}{2}\|\boldsymbol{\beta} - \boldsymbol{\beta}^{t-1}\|_2^2 + \lambda_\beta(\|\boldsymbol{\beta}\|_1 - \|\boldsymbol{\beta}^{t-1}\|_1) \\ &\leq F(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1}) + \langle \boldsymbol{\beta} - \boldsymbol{\beta}^{t-1}, \partial\mathcal{L}(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1})/\partial\boldsymbol{\beta} \rangle \\ &\quad + \frac{1}{2\eta}\|\boldsymbol{\beta} - \boldsymbol{\beta}^{t-1}\|_2^2 + \lambda_\beta(\|\boldsymbol{\beta}\|_1 - \|\boldsymbol{\beta}^{t-1}\|_1) \end{aligned} \quad (32)$$

for any  $\boldsymbol{\beta}$ , where we have used  $\sigma_\beta \leq 1/\eta$ .

Furthermore, for any given  $b \geq 0$ , let  $\boldsymbol{\beta}^t = b\tilde{\boldsymbol{\beta}}^{t-1} + (1-b)\boldsymbol{\beta}^{t-1}$ , where recall that

$$\tilde{\boldsymbol{\beta}}^{t-1} \equiv \operatorname{argmin}_\beta \langle \partial\mathcal{L}(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1})/\partial\boldsymbol{\beta}, \boldsymbol{\beta} \rangle + \lambda_\beta\|\boldsymbol{\beta}\|_1, \|\boldsymbol{\beta}\|_1 \leq Q^{t-1} \quad (33)$$

with  $Q^{t-1} = F(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1})/\lambda_\beta$ , it holds that

$$\begin{aligned} F(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^t) &\leq F(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1}) + b\langle \partial\mathcal{L}(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1})/\partial\boldsymbol{\beta}, \tilde{\boldsymbol{\beta}}^{t-1} - \boldsymbol{\beta}^{t-1} \rangle + \frac{b^2}{2\eta}\|\tilde{\boldsymbol{\beta}}^{t-1} - \boldsymbol{\beta}^{t-1}\|_2^2 \\ &\quad + \lambda_\beta(\|b\tilde{\boldsymbol{\beta}}^{t-1} + (1-b)\boldsymbol{\beta}^{t-1}\|_1 - \|\boldsymbol{\beta}^{t-1}\|_1) \\ &\leq F(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1}) + b\langle \partial\mathcal{L}(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1})/\partial\boldsymbol{\beta}, \tilde{\boldsymbol{\beta}}^{t-1} - \boldsymbol{\beta}^{t-1} \rangle + \frac{b^2}{2\eta}\|\tilde{\boldsymbol{\beta}}^{t-1} - \boldsymbol{\beta}^{t-1}\|_2^2 \\ &\quad + b\lambda_\beta(\|\tilde{\boldsymbol{\beta}}^{t-1} - \boldsymbol{\beta}^{t-1}\|_1). \end{aligned} \quad (34)$$

To minimize the right hand side, we set

$$\begin{aligned} b &\equiv \frac{\eta\{-\langle \partial\mathcal{L}(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1})/\partial\boldsymbol{\beta}, \tilde{\boldsymbol{\beta}}^{t-1} - \boldsymbol{\beta}^{t-1} \rangle - \lambda_\beta(\|\tilde{\boldsymbol{\beta}}^{t-1} - \boldsymbol{\beta}^{t-1}\|_1)\}}{\|\tilde{\boldsymbol{\beta}}^{t-1} - \boldsymbol{\beta}^{t-1}\|_2^2} \\ &= \frac{\eta g_\beta(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1}, Q^{t-1})}{\|\tilde{\boldsymbol{\beta}}^{t-1} - \boldsymbol{\beta}^{t-1}\|_2^2}. \end{aligned}$$

Note that  $g_\beta(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1}, Q^{t-1}) \geq 0$  due to its definition and the definition of  $\tilde{\boldsymbol{\beta}}^{t-1}$ , hence  $b \geq 0$ .

Note that  $F(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1}) = \mathcal{L}(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1}) - L + \lambda_\Theta\|\boldsymbol{\Theta}^{t-1}\|_* + \lambda_\beta\|\boldsymbol{\beta}^{t-1}\|_1 \geq \lambda_\beta\|\boldsymbol{\beta}^{t-1}\|_1$ , hence  $\|\boldsymbol{\beta}^{t-1}\|_1 \leq Q^{t-1}$ . Since  $\tilde{\boldsymbol{\beta}}^{t-1}$  is the minimizer of (33), we get

$$\begin{aligned} &g_\beta(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1}, Q^{t-1}) \\ &= \langle \partial\mathcal{L}(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1})/\partial\boldsymbol{\beta}, \boldsymbol{\beta}^{t-1} - \tilde{\boldsymbol{\beta}}^{t-1} \rangle + \lambda_\beta\|\boldsymbol{\beta}^{t-1} - \tilde{\boldsymbol{\beta}}^{t-1}\|_1 \\ &= \langle \partial\mathcal{L}(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1})/\partial\boldsymbol{\beta}, \boldsymbol{\beta}^{t-1} - \tilde{\boldsymbol{\beta}}^{t-1} \rangle + \lambda_\beta\|\boldsymbol{\beta}^{t-1}\|_1 - \lambda_\beta\|\tilde{\boldsymbol{\beta}}^{t-1}\|_1 \geq 0. \end{aligned}$$

Plug  $b$  in (34), we have

$$F(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^t) - F(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1}) \leq -\frac{\eta g_\beta(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1}, Q^{t-1})^2}{2\|\tilde{\boldsymbol{\beta}}^{t-1} - \boldsymbol{\beta}^{t-1}\|_2^2} \leq -\frac{\eta g_\beta(\boldsymbol{\Theta}^{t-1}, \boldsymbol{\beta}^{t-1}, Q^{t-1})^2}{2(2Q^{t-1})^2}. \quad (35)$$

The last equality holds because  $\|\boldsymbol{\beta}^{t-1}\|_1 \leq Q^{t-1}$  and  $\|\tilde{\boldsymbol{\beta}}^{t-1}\|_1 \leq Q^{t-1}$ , and hence  $\|\tilde{\boldsymbol{\beta}}^{t-1} - \boldsymbol{\beta}^{t-1}\|_2^2 \leq (2Q^{t-1})^2$ . This proves the result.

**Lemma A.23** *Assume*

$$\text{vec}(\Theta)^T E\{\mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\} \text{vec}(\Theta) > 4\sqrt{(d_{H_2} + d_{S_2}^2)\log(mn)/(mn)} \|\Theta\|_F^2$$

for all  $\Theta$  that satisfies  $\|\Theta\|_{\max} \leq 2a$  and  $\tilde{\beta}, \tilde{\Theta}$  that satisfy  $\|\tilde{\beta}\|_\infty \leq a$  and  $\|\tilde{\Theta}\|_{\max} \leq a$ . Recall that  $\sigma_\Theta = 2d_{H_2} + 2d_{S_2}^2$ ,  $R^{t-1} = F(\Theta^{t-1}, \beta^t)/\lambda_\Theta$ , and

$$g_\Theta(\Theta^{t-1}, \beta^t, R^{t-1}) = \langle \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial\Theta, \Theta^{t-1} - \tilde{\Theta}^{t-1} \rangle + \lambda_\Theta \|\Theta^{t-1} - \tilde{\Theta}^{t-1}\|_*.$$

where

$$\tilde{\Theta}^{t-1} \equiv \text{argmin}_\Theta \langle \partial\mathcal{L}(\Theta^{t-1}, \beta^t) / \partial\Theta, \Theta \rangle + \lambda_\Theta \|\Theta\|_*, \|\Theta\|_* \leq R^{t-1}.$$

Select  $0 < \eta_1 < 1/\sigma_\Theta$ . Then we have

$$F(\Theta^t, \beta^t) \leq F(\Theta^{t-1}, \beta^t) - \frac{\eta_1 g_\Theta(\Theta^{t-1}, \beta^t, R^{t-1})^2}{2(2R^{t-1})^2},$$

with probability at least  $1 - 2(mn)^{-1}$ .

Proof: Because

$$E\{\text{vec}(\Theta)^T \mathbf{F}_\Theta(\mathbf{X}, \mathbf{Y}|\tilde{\Theta}, \tilde{\beta})\} \text{vec}(\Theta) > 4\sqrt{(d_{H_2} + d_{S_2}^2)\log(mn)/(mn)} \|\Theta\|_F^2,$$

then the second derivative of  $\mathcal{L}(\Theta, \beta)$  with respect to  $\text{vec}(\Theta)$  is positive definite with probability at least  $1 - 2(mn)^{-1}$  by Lemma A.21, and hence we majorize  $F(\Theta, \beta^t)$  at  $\Theta^{t-1}$  as

$$\begin{aligned} F(\Theta, \beta^t) &\leq F(\Theta^{t-1}, \beta^t) + \langle \Theta - \Theta^{t-1}, \partial\mathcal{L}(\Theta^{t-1}, \beta^t) / \partial\Theta \rangle \\ &\quad + \frac{\sigma_\Theta}{2} \|\Theta - \Theta^{t-1}\|_F^2 + \lambda_\Theta (\|\Theta\|_* - \|\Theta^{t-1}\|_*) \\ &\leq F(\Theta^{t-1}, \beta^t) + \langle \Theta - \Theta^{t-1}, \partial\mathcal{L}(\Theta^{t-1}, \beta^t) / \partial\Theta \rangle \\ &\quad + \frac{1}{2\eta_1} \|\Theta - \Theta^{t-1}\|_F^2 + \lambda_\Theta (\|\Theta\|_* - \|\Theta^{t-1}\|_*) \end{aligned}$$

for any  $\Theta$ .

Furthermore, for any given  $b \geq 0$ , let  $\Theta^t = b\tilde{\Theta}^{t-1} + (1-b)\Theta^{t-1}$ , where recall that

$$\tilde{\Theta}^{t-1} \equiv \text{argmin}_\Theta \langle \partial\mathcal{L}(\Theta^{t-1}, \beta^t) / \partial\Theta, \Theta \rangle + \lambda_\Theta \|\Theta\|_*, \|\Theta\|_* \leq R^{t-1} \quad (36)$$

with  $R^{t-1} = F(\Theta^{t-1}, \beta^t)/\lambda_\Theta$ , it holds that

$$\begin{aligned} F(\Theta^t, \beta^t) &\leq F(\Theta^{t-1}, \beta^t) + b \langle \partial\mathcal{L}(\Theta^{t-1}, \beta^t) / \partial\Theta, \tilde{\Theta}^{t-1} - \Theta^{t-1} \rangle + \frac{b^2}{2\eta_1} \|\tilde{\Theta}^{t-1} - \Theta^{t-1}\|_F^2 \\ &\quad + \lambda_\Theta (\|b\tilde{\Theta}^{t-1} + (1-b)\Theta^{t-1}\|_* - \|\Theta^{t-1}\|_*) \\ &\leq F(\Theta^{t-1}, \beta^t) + b \langle \partial\mathcal{L}(\Theta^{t-1}, \beta^t) / \partial\Theta, \tilde{\Theta}^{t-1} - \Theta^{t-1} \rangle + \frac{b^2}{2\eta_1} \|\tilde{\Theta}^{t-1} - \Theta^{t-1}\|_F^2 \\ &\quad + b\lambda_\Theta (\|\tilde{\Theta}^{t-1} - \Theta^{t-1}\|_*). \end{aligned} \quad (37)$$

To minimize the right hand side, we set

$$b \equiv \frac{\eta_1 \{-\langle \partial\mathcal{L}(\Theta^{t-1}, \beta^t) / \partial\Theta, \tilde{\Theta}^{t-1} - \Theta^{t-1} \rangle - \lambda_\Theta (\|\tilde{\Theta}^{t-1} - \Theta^{t-1}\|_*)\}}{\|\tilde{\Theta}^{t-1} - \Theta^{t-1}\|_F^2}$$

$$= \frac{\eta_1 g_{\Theta}(\Theta^{t-1}, \beta^t, R^{t-1})}{\|\tilde{\Theta}^{t-1} - \Theta^{t-1}\|_F^2}.$$

Note that  $g_{\Theta}(\Theta^{t-1}, \beta^t, Q^{t-1}) \geq 0$  due to its definition and the definition of  $\tilde{\Theta}^{t-1}$ , hence  $b \geq 0$ . Note that  $F(\Theta^{t-1}, \beta^t) = \mathcal{L}(\Theta^{t-1}, \beta^t) - L + \lambda_{\Theta} \|\Theta^{t-1}\|_* + \lambda_{\beta} \|\beta^t\|_1 \geq \lambda_{\Theta} \|\Theta^{t-1}\|_*$ , hence  $\|\Theta^{t-1}\|_* \leq R^{t-1}$ . Since  $\tilde{\Theta}^{t-1}$  is the minimizer of (36), we get

$$\begin{aligned} & g_{\Theta}(\Theta^{t-1}, \beta^t, R^{t-1}) \\ \equiv & \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \Theta, \Theta^{t-1} - \tilde{\Theta}^{t-1} \rangle + \lambda_{\Theta} \|\Theta^{t-1} - \tilde{\Theta}^{t-1}\|_* \\ \geq & \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \Theta, \Theta^{t-1} - \tilde{\Theta}^{t-1} \rangle + \lambda_{\Theta} \|\Theta^{t-1}\|_* - \lambda_{\Theta} \|\tilde{\Theta}^{t-1}\|_* \geq 0. \end{aligned}$$

Plug  $b$  in (37), we have

$$F(\Theta^t, \beta^t) - F(\Theta^{t-1}, \beta^t) \leq -\frac{\eta_1}{2} \frac{g_{\Theta}(\Theta^{t-1}, \beta^t, R^{t-1})^2}{\|\tilde{\Theta}^{t-1} - \Theta^{t-1}\|_F^2} \leq -\frac{\eta_1}{2} \frac{g_{\Theta}(\Theta^{t-1}, \beta^t, R^{t-1})^2}{(2R^{t-1})^2}. \quad (38)$$

The last equality holds because  $\|\Theta^{t-1}\|_1 \leq R^{t-1}$  and  $\|\tilde{\Theta}^{t-1}\|_1 \leq R^{t-1}$ , and hence  $\|\tilde{\Theta}^{t-1} - \Theta^{t-1}\|_2^2 \leq (2R^{t-1})^2$ . This proves the result.

## Appendix H. Proof of Theorem 3

Proof: Let

$$\tilde{\beta}^t = \operatorname{argmin} \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \beta, \beta \rangle + \lambda_{\beta} \|\beta\|_1, \|\beta\|_1 \leq Q^t,$$

and recall that

$$\tilde{\Theta}^{t-1} \equiv \operatorname{argmin}_{\Theta} \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \Theta, \Theta \rangle + \lambda_{\Theta} \|\Theta\|_*, \|\Theta\|_* \leq R^{t-1}$$

as defined in (36). Because

$$\begin{aligned} & E\{\{\beta^T, \operatorname{vec}(\Theta)^T\} \partial^2 \mathcal{L}(\Theta, \beta) / \partial [\{\operatorname{vec}(\Theta)^T, \beta^T\}^T]^{\otimes 2} \{\beta^T, \operatorname{vec}(\Theta)^T\}^T\} \\ \geq & 4(c_0 + 1) \sqrt{(d_{H_2} + d_{S_2}^2) \log\{\max(p, mn)\} / (mn)} \|\{\beta^T, \operatorname{vec}(\Theta)^T\}^T\|_2^2, \end{aligned}$$

we have

$$\begin{aligned} & F(\hat{\Theta}, \hat{\beta}) - F(\Theta^{t-1}, \beta^t) \\ \geq & \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \Theta, \hat{\Theta} - \Theta^{t-1} \rangle + \lambda_{\Theta} (\|\hat{\Theta}\|_* - \|\Theta^{t-1}\|_*) \\ & + \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \beta, \hat{\beta} - \beta^t \rangle + \lambda_{\beta} (\|\hat{\beta}\|_1 - \|\beta^t\|_1), \end{aligned}$$

with probability at least  $1 - 2(mn + p)^{-1}$  by Lemma A.21. Hence

$$\begin{aligned} F(\Theta^{t-1}, \beta^t) - F(\hat{\Theta}, \hat{\beta}) & \leq \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \Theta, \Theta^{t-1} - \hat{\Theta} \rangle + \lambda_{\Theta} (\|\Theta^{t-1}\|_* - \|\hat{\Theta}\|_*) \\ & \quad + \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \beta, \beta^t - \hat{\beta} \rangle + \lambda_{\beta} (\|\beta^t\|_1 - \|\hat{\beta}\|_1) \\ & \leq \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \Theta, \Theta^{t-1} - \tilde{\Theta}^{t-1} \rangle + \lambda_{\Theta} (\|\Theta^{t-1}\|_* - \|\tilde{\Theta}^{t-1}\|_*) \\ & \quad + \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \beta, \beta^t - \tilde{\beta}^t \rangle + \lambda_{\beta} (\|\beta^t\|_1 - \|\tilde{\beta}^t\|_1) \\ & \leq \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \Theta, \Theta^{t-1} - \tilde{\Theta}^{t-1} \rangle + \lambda_{\Theta} (\|\Theta^{t-1} - \tilde{\Theta}^{t-1}\|_*) \end{aligned}$$



$$\begin{aligned}
 & + \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \beta, \beta^t - \tilde{\beta}^t \rangle + \lambda_\beta (\|\beta^t - \tilde{\beta}^t\|_1) \\
 = & g_\Theta(\Theta^{t-1}, \beta^t, R^{t-1}) + g_\beta(\Theta^{t-1}, \beta^t, Q^t),
 \end{aligned} \tag{39}$$

where recall that  $g_\beta(\Theta^{t-1}, \beta^t, Q^t) = \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \beta, \beta^t - \tilde{\beta}^t \rangle + \lambda_\beta (\|\beta^t - \tilde{\beta}^t\|_1)$  and  $g_\Theta(\Theta^{t-1}, \beta^t, R^{t-1}) = \langle \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \Theta, \Theta^{t-1} - \tilde{\Theta}^{t-1} \rangle + \lambda_\Theta \|\Theta^{t-1} - \tilde{\Theta}^{t-1}\|_*$ . The second inequality holds because  $F(\Theta, \beta)$  is minimized at  $\hat{\Theta}$  and  $\hat{\beta}$ , and hence  $\|\hat{\beta}\|_1 \leq F(\hat{\Theta}, \hat{\beta}) / \lambda_\beta \leq F(\Theta^t, \beta^t) / \lambda_\beta = Q^t$  and  $\|\hat{\Theta}\|_* \leq F(\hat{\Theta}, \hat{\beta}) / \lambda_\Theta \leq F(\Theta^{t-1}, \beta^t) / \lambda_\Theta \leq R^{t-1}$ . Furthermore,

$$\begin{aligned}
 g_\beta(\Theta^{t-1}, \beta^t, Q^t) & = \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \beta, \beta^t - \tilde{\beta}^t \rangle + \lambda_\beta (\|\beta^t - \tilde{\beta}^t\|_1) \\
 & = \langle \partial \mathcal{L}(\Theta^t, \beta^t) / \partial \beta, \beta^t - \tilde{\beta}^t \rangle + \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \beta - \partial \mathcal{L}(\Theta^t, \beta^t) / \partial \beta, \beta^t - \tilde{\beta}^t \rangle \\
 & \quad + \lambda_\beta (\|\beta^t - \tilde{\beta}^t\|_1) \\
 & \leq \langle \partial \mathcal{L}(\Theta^t, \beta^t) / \partial \beta, \beta^t - \tilde{\beta}^t \rangle + \lambda_\beta \|\beta^t - \tilde{\beta}^t\|_1 \\
 & \quad + \|\partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \beta - \partial \mathcal{L}(\Theta^t, \beta^t) / \partial \beta\|_2 \|\beta^t - \tilde{\beta}^t\|_2 \\
 & \leq g_\beta(\Theta^t, \beta^t, Q^t) + 2\sigma_{\beta\Theta} Q^t \|\Theta^{t-1} - \Theta^t\|_F.
 \end{aligned} \tag{40}$$

The last inequality holds by Remark 1 and  $\|\beta^t - \tilde{\beta}^t\|_2 \leq \|\beta^t - \tilde{\beta}^t\|_1 \leq 2Q^t$ . Furthermore, because  $\Theta^t$  is the minimizer of  $\frac{1}{2} \|\Theta - \Theta^{t-1}\|_F^2 + \eta_1 \frac{\partial \mathcal{L}(\Theta^{t-1}, \beta^t)}{\partial \Theta} \|\Theta\|_*^2 + \eta_1 \lambda_\Theta \|\Theta\|_*$ , we have

$$\Theta^t - \Theta^{t-1} + \eta_1 \frac{\partial \mathcal{L}(\Theta^{t-1}, \beta^t)}{\partial \Theta} + \eta_1 \lambda_\Theta \partial \|\Theta^t\|_* / \partial \Theta = 0.$$

Therefore,

$$\begin{aligned}
 & F(\Theta^{t-1}, \beta^t) - F(\Theta^t, \beta^t) \\
 = & \langle \partial \mathcal{L}(\Theta^t, \beta^t) / \partial \Theta + \lambda_\Theta \partial \|\Theta^t\|_* / \partial \Theta, \Theta^{t-1} - \Theta^t \rangle \\
 & + \{\text{vec}(\Theta^{t-1}) - \text{vec}(\Theta^t)\}^T \frac{\partial^2 \mathcal{L}(\Theta^*, \beta^t)}{2 \partial \text{vec}(\Theta) \text{vec}(\Theta)^T} \{\text{vec}(\Theta^{t-1}) - \text{vec}(\Theta^t)\} \\
 = & \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \Theta + \lambda_\Theta \partial \|\Theta^t\|_* / \partial \Theta, \Theta^{t-1} - \Theta^t \rangle \\
 & + \langle \partial \mathcal{L}(\Theta^t, \beta^t) / \partial \Theta - \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \Theta, \Theta^{t-1} - \Theta^t \rangle \\
 & + \{\text{vec}(\Theta^{t-1}) - \text{vec}(\Theta^t)\}^T \frac{\partial^2 \mathcal{L}(\Theta^*, \beta^t)}{2 \partial \text{vec}(\Theta) \text{vec}(\Theta)^T} \{\text{vec}(\Theta^{t-1}) - \text{vec}(\Theta^t)\} \\
 = & \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \Theta + \lambda_\Theta \partial \|\Theta^t\|_* / \partial \Theta, \Theta^{t-1} - \Theta^t \rangle \\
 & - \{\text{vec}(\Theta^{t-1}) - \text{vec}(\Theta^t)\}^T \frac{\partial^2 \mathcal{L}(\Theta^{**}, \beta^t)}{\partial \text{vec}(\Theta) \text{vec}(\Theta)^T} \{\text{vec}(\Theta^{t-1}) - \text{vec}(\Theta^t)\} \\
 & + \{\text{vec}(\Theta^{t-1}) - \text{vec}(\Theta^t)\}^T \frac{\partial^2 \mathcal{L}(\Theta^*, \beta^t)}{2 \partial \text{vec}(\Theta) \text{vec}(\Theta)^T} \{\text{vec}(\Theta^{t-1}) - \text{vec}(\Theta^t)\} \\
 \geq & \langle \partial \mathcal{L}(\Theta^{t-1}, \beta^t) / \partial \Theta + \lambda_\Theta \partial \|\Theta^t\|_* / \partial \Theta, \Theta^{t-1} - \Theta^t \rangle \\
 & - \{\text{vec}(\Theta^{t-1}) - \text{vec}(\Theta^t)\}^T \frac{\partial^2 \mathcal{L}(\Theta^{**}, \beta^t)}{\partial \text{vec}(\Theta) \text{vec}(\Theta)^T} \{\text{vec}(\Theta^{t-1}) - \text{vec}(\Theta^t)\} \\
 = & 1/\eta_1 \|\Theta^{t-1} - \Theta^t\|_F^2 - \{\text{vec}(\Theta^{t-1}) - \text{vec}(\Theta^t)\}^T \frac{\partial^2 \mathcal{L}(\Theta^{**}, \beta^t)}{\partial \text{vec}(\Theta) \partial \text{vec}(\Theta)^T} \{\text{vec}(\Theta^{t-1}) - \text{vec}(\Theta^t)\} \\
 \geq & (1/\eta_1 - \sigma_\Theta) \|\Theta^{t-1} - \Theta^t\|_F^2,
 \end{aligned} \tag{41}$$

where  $\Theta^*$  and  $\Theta^{**}$  are points on the line between  $\Theta^{t-1}$  and  $\Theta^t$ . Hence

$$\|\Theta^{t-1} - \Theta^t\|_F^2 \leq \{F(\Theta^{t-1}, \beta^t) - F(\Theta^t, \beta^t)\} / (1/\eta_1 - \sigma_\Theta).$$

The last inequality holds by the fact that  $\|\partial^2 \mathcal{L}(\Theta^{**}, \beta^t) / \partial \text{vec}(\Theta) \partial \text{vec}(\Theta)^T\|_2 \leq \sigma_{\Theta}$  as shown in Remark 1. Combining (40), (41), Lemmas A.22 and A.23, we obtain that with probability at least  $1 - 2 \max(p, mn)^{-1} - 2(mn)^{-1}$ ,

$$\begin{aligned}
 & \{g_{\Theta}(\Theta^{t-1}, \beta^t, R^{t-1}) + g_{\beta}(\Theta^{t-1}, \beta^t, Q^t)\}^2 \\
 \leq & 4\{g_{\Theta}(\Theta^{t-1}, \beta^t, R^{t-1})^2 + g_{\beta}(\Theta^t, \beta^t, Q^t)^2 + 4\sigma_{\beta_{\Theta}}^2(Q^t)^2\|\Theta^{t-1} - \Theta^t\|_F^2\} \\
 \leq & 4[2(2R^{t-1})^2/\eta_1\{F(\Theta^{t-1}, \beta^t) - F(\Theta^t, \beta^t)\} \\
 & + 2(2Q^t)^2/\eta\{F(\Theta^t, \beta^t) - F(\Theta^t, \beta^{t+1})\} \\
 & + 4\sigma_{\beta_{\Theta}}^2(Q^t)^2\{F(\Theta^{t-1}, \beta^t) - F(\Theta^t, \beta^t)\}/(1/\eta_1 - \sigma_{\Theta})] \\
 \leq & C(t)\{F(\Theta^{t-1}, \beta^t) - F(\Theta^t, \beta^t)\} + C(t)\{F(\Theta^t, \beta^t) - F(\Theta^t, \beta^{t+1})\} \\
 = & C(t)\{F(\Theta^{t-1}, \beta^t) - F(\Theta^t, \beta^{t+1})\},
 \end{aligned}$$

where  $C(t) = \max\{32(R^{t-1})^2/\eta_1 + 16\sigma_{\beta_{\Theta}}^2(Q^t)^2/(1/\eta_1 - \sigma_{\Theta}), 32(Q^t)^2/\eta\}$ . Combining with (39), we have

$$\{F(\Theta^{t-1}, \beta^t) - F(\widehat{\Theta}, \widehat{\beta})\}^2 \leq C(t)\{F(\Theta^{t-1}, \beta^t) - F(\widehat{\Theta}, \widehat{\beta}) - F(\Theta^t, \beta^{t+1}) + F(\widehat{\Theta}, \widehat{\beta})\}$$

with probability at least  $1 - 2 \max(p, mn)^{-1} - 2(mn)^{-1} - 2(mn + p)^{-1}$ . Let  $\Delta_t = F(\Theta^{t-1}, \beta^t) - F(\widehat{\Theta}, \widehat{\beta})$ , the above inequality can be written as

$$\Delta_{t+1} \leq \Delta_t - \frac{1}{C(t)}(\Delta_t)^2 = \Delta_t \left\{1 - \frac{\Delta_t}{C(t)}\right\} \leq \Delta_t \left\{1 + \frac{\Delta_t}{C(t)}\right\}^{-1}.$$

Taking inverse, we get  $\Delta_{t+1}^{-1} \geq \Delta_t^{-1} + C(t)^{-1}$ , hence  $\Delta_{t+1}^{-1} \geq \Delta_0^{-1} + \sum_{k=0}^t C(k)^{-1}$ , which leads to

$$\Delta_{t+1} \leq \frac{1}{1/\{F(\Theta^0, \beta^0) - F(\widehat{\Theta}, \widehat{\beta})\} + \sum_{k=0}^t 1/C(k)}.$$

Hence when

$$\sum_{t=0}^{T-1} 1/C(t) \geq \left\{\frac{1}{\epsilon} - \frac{1}{F(\Theta^0, \beta^0) - F(\widehat{\Theta}, \widehat{\beta})}\right\},$$

we have

$$F(\Theta^T, \beta^T) - F(\widehat{\Theta}, \widehat{\beta}) \leq \epsilon,$$

with probability at least  $1 - 2 \max(p, mn)^{-1} - 2(mn)^{-1} - 2(mn + p)^{-1} \geq 1 - 6(mn + p)^{-1}$ . Hence  $F(\Theta^T, \beta^T)$  is the  $\epsilon$ -optimal solution for (1) with probability at least  $1 - 6(mn + p)^{-1}$ . This proves the result.

## Appendix I. Proof of the Consistency of Pseudo-likelihood Estimator in Regression with Finite Dimensional Parameter

**Theorem A.1** *Assume Condition (C1)–(C6) hold and  $\beta$  is finite dimensional. Furthermore, assume the conditional density of  $Y_i$  given  $\mathbf{X}_i$  is  $f(Y_i, \beta_0^T \mathbf{X}_i)$ . Let  $\widehat{\beta}$  be the maximizer for the pseudo-likelihood*

$$\mathcal{L}(\beta) \equiv n^{-1} \sum_{i=1}^n \mathcal{L}_i(\beta),$$

where  $\mathcal{L}_i(\beta) = R_i \ell_i(\beta)$ ,

$$\ell_i(\beta) = \left[ \log\{f(Y_i, \beta^T \mathbf{X}_i)\} - \log \left\{ \int f(Y_i, \beta^T \mathbf{X}) g(\mathbf{X}) d\mathbf{X} \right\} \right],$$

$R_i$  is the missing indicator for the  $i$ th observation. Assume  $\mathbf{X}_i$  and  $R_i$  are independent given  $Y_i$ . We have  $\hat{\boldsymbol{\beta}} \xrightarrow{p} \boldsymbol{\beta}_0$ .

Proof: Everywhere in this proof,  $E$  denotes an expectation taken with respect to the original true distribution that governs the data generation process  $f(Y, \boldsymbol{\beta}_0^T \mathbf{X})g(\mathbf{X})$ , for example, for any function  $\mathbf{a}(\mathbf{X}, Y)$ ,  $E\{\mathbf{a}(\mathbf{X}, Y)\} = \int \mathbf{a}(\mathbf{x}, y)f(y, \boldsymbol{\beta}_0^T \mathbf{x})g(\mathbf{x})d\mu(\mathbf{x}, y)$ . Let  $f_2(Y_i, \boldsymbol{\beta}^T \mathbf{X}_i)$  be the derivative of  $f(Y_i, \boldsymbol{\beta}^T \mathbf{X}_i)$  with respect to  $\boldsymbol{\beta}^T \mathbf{X}_i$ ,  $S_2(Y_i, \mathbf{X}_i, \boldsymbol{\beta})$  be the derivative of  $\log f(Y_i, \boldsymbol{\beta}^T \mathbf{X}_i)$  with respect to  $\boldsymbol{\beta}^T \mathbf{X}_i$ , define

$$\begin{aligned} \mathbf{S}(Y_i, \mathbf{X}_i, \boldsymbol{\beta}) &\equiv \frac{\partial \mathcal{L}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \\ &= n^{-1} \sum_{i=1}^n R_i \left\{ \frac{f_2(Y_i, \boldsymbol{\beta}^T \mathbf{X}_i)}{f(Y_i, \boldsymbol{\beta}^T \mathbf{X}_i)} \mathbf{X}_i - \frac{\int f_2(Y_i, \boldsymbol{\beta}^T \mathbf{X})g(\mathbf{X})\mathbf{X}d\mathbf{X}}{\int f(Y_i, \boldsymbol{\beta}^T \mathbf{X})g(\mathbf{X})d\mathbf{X}} \right\} \\ &= n^{-1} \sum_{i=1}^n R_i S_2(Y_i, \mathbf{X}_i, \boldsymbol{\beta}) \mathbf{X}_i - R_i E_{\boldsymbol{\beta}} \{S_2(Y_i, \mathbf{X}, \boldsymbol{\beta}) \mathbf{X} | Y_i\}. \end{aligned}$$

Here, the notation  $E_{\boldsymbol{\beta}}$  means an expectation is taken with respect to the distribution  $f(Y, \boldsymbol{\beta}^T \mathbf{X})g(\mathbf{X})$ . It is clear that

$$\frac{\partial E\{\mathcal{L}(\boldsymbol{\beta})\}}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} = E \left\{ \frac{\partial \mathcal{L}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right\} = E\{\mathbf{S}(Y_i, \mathbf{X}_i, \boldsymbol{\beta}_0)\} = 0. \quad (42)$$

Also,

$$\begin{aligned} \frac{\partial^2 \mathcal{L}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} &= n^{-1} \sum_{i=1}^n R_i \left( \mathbf{H}(Y_i, \mathbf{X}_i, \boldsymbol{\beta}) - E_{\boldsymbol{\beta}} \{ \mathbf{H}(Y_i, \mathbf{X}_i, \boldsymbol{\beta}) | Y_i \} \right. \\ &\quad \left. - E_{\boldsymbol{\beta}} \{ \mathbf{S}(Y_i, \mathbf{X}_i, \boldsymbol{\beta})^{\otimes 2} | Y_i \} + [E_{\boldsymbol{\beta}} \{ \mathbf{S}(Y_i, \mathbf{X}_i, \boldsymbol{\beta}) | Y_i \}]^{\otimes 2} \right), \end{aligned} \quad (43)$$

where  $\mathbf{H}(Y_i, \mathbf{X}_i, \boldsymbol{\beta})$  is the second partial derivative of  $\log\{f(Y_i, \boldsymbol{\beta}^T \mathbf{X}_i)\}$  with respect to  $\boldsymbol{\beta}$ . Now using the independence between  $\mathbf{X}_i$  and  $R_i$  given  $Y_i$ , we have

$$\begin{aligned} E \left\{ \frac{\partial^2 \mathcal{L}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right\} &= E \left( E[R_i \mathbf{H}(Y_i, \mathbf{X}_i, \boldsymbol{\beta}) | Y_i] - E(R_i | Y_i) E_{\boldsymbol{\beta}_0} \{ \mathbf{H}(Y_i, \mathbf{X}_i, \boldsymbol{\beta}_0) | Y_i \} \right. \\ &\quad \left. - E \left\{ E(R_i | Y_i) \left( E \{ \mathbf{S}(Y_i, \mathbf{X}_i, \boldsymbol{\beta}_0)^{\otimes 2} | Y_i \} - [E\{ \mathbf{S}(Y_i, \mathbf{X}_i, \boldsymbol{\beta}_0) | Y_i \}]^{\otimes 2} \right) \right\} \right) \\ &= -E \left\{ E(R_i | Y_i) \left( E \{ \mathbf{S}(Y_i, \mathbf{X}_i, \boldsymbol{\beta}_0)^{\otimes 2} | Y_i \} - [E\{ \mathbf{S}(Y_i, \mathbf{X}_i, \boldsymbol{\beta}_0) | Y_i \}]^{\otimes 2} \right) \right\} \\ &= -E \left[ E(R_i | Y_i) \text{var} \{ \mathbf{S}(Y_i, \mathbf{X}_i, \boldsymbol{\beta}_0) | Y_i \} \right], \end{aligned} \quad (44)$$

which is negative definite. Note that  $E_{\boldsymbol{\beta}_0} = E$  according to our definition of the two notations. Now by Taylor expansion, let  $\mathbf{r}$  be a vector with  $\|\mathbf{r}\|_2 = 1$ , assume  $\mathbf{r}$  times the third derivative of  $\mathcal{L}(\boldsymbol{\beta})$  with respect to  $\boldsymbol{\beta}$  has bounded  $L_2$  norm, for any  $\epsilon \in (0, 1/2)$ , we have

$$\begin{aligned} &\mathcal{L}(\boldsymbol{\beta}_0 + \mathbf{r}n^{-1/2+\epsilon}) \\ &= \mathcal{L}(\boldsymbol{\beta}_0) + \frac{\partial \mathcal{L}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \mathbf{r}n^{-1/2+\epsilon} + \frac{1}{2} \mathbf{r}^T \frac{\partial^2 \mathcal{L}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \mathbf{r}n^{-1+2\epsilon} + O_p(n^{-3/2+3\epsilon}) \\ &= \mathcal{L}(\boldsymbol{\beta}_0) + \left[ E \left\{ \frac{\partial \mathcal{L}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right\} + O_p(n^{-1/2}) \right] \mathbf{r}n^{-1/2+\epsilon} \\ &+ \frac{1}{2} \mathbf{r}^T \left[ E \left\{ \frac{\partial^2 \mathcal{L}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right\} + o_p(1) \right] \mathbf{r}n^{-1+2\epsilon} + O_p(n^{-3/2+3\epsilon}) \end{aligned}$$

$$= \mathcal{L}(\beta_0) + \frac{1}{2} \mathbf{r}^T E \left\{ \frac{\partial^2 \mathcal{L}(\beta)}{\partial \beta \partial \beta^T} \Big|_{\beta=\beta_0} \right\} \mathbf{r} n^{-1+2\epsilon} + o_p(n^{-1+2\epsilon}), \quad (45)$$

where the second equality holds by the central limit theorem, the third equality holds by (42). Now by (44),  $E \left\{ \frac{\partial^2 \mathcal{L}(\beta)}{\partial \beta \partial \beta^T} \Big|_{\beta=\beta_0} \right\}$  is negative definite. Hence,  $\mathcal{L}(\beta_0) > \mathcal{L}(\beta_0 + \mathbf{r}/\sqrt{n})$  in probability when  $n$  is sufficiently large for any  $\mathbf{r}$  with  $\|\mathbf{r}\|_2 = 1$ . Therefore, there is a maximizer for  $\mathcal{L}(\beta)$  in the ball with center  $\beta_0$  and radius  $n^{-/2+\epsilon}$ . Let the maximizer be  $\hat{\beta}$ . Obviously  $\|\hat{\beta} - \beta_0\| \leq n^{-/2+\epsilon}$  in probability, hence  $\hat{\beta} \xrightarrow{p} \beta_0$ .

## References

- J. Abernethy, F. Bach, T. Evgeniou, and J.-P. Vert. A new approach to collaborative filtering: Operator estimation with spectral regularization. *Journal of Machine Learning Research*, 10:803–826, 2009.
- Afonso S Bandeira, Ramon Van Handel, et al. Sharp nonasymptotic bounds on the norm of random matrices with independent entries. *The Annals of Probability*, 44(4):2479–2506, 2016.
- X. Bi, A. Qu, J. Wang, and X. Shen. A group-specific recommender system. *Journal of the American Statistical Association*, 112:1344–1353, 2016.
- P. J. Bickel, Y. Ritov, A. B. Tsybakov, et al. Simultaneous analysis of lasso and dantzig selector. *The Annals of Statistics*, 37(4):1705–1732, 2009.
- T. Cai, T. T. Cai, and A. Zhang. Structured matrix completion with applications to genomic data integration. *Journal of the American Statistical Association*, 111:621–633, 2016.
- T. T. Cai and W.-X. Zhou. Matrix completion via max-norm constrained optimization. *Electronic Journal of Statistics*, 10:1493–1525, 2016.
- E. J. Candés and Y. Plan. Matrix completion with noise. *Proceedings of the IEEE*, 98:925–936, 2010.
- E. J. Candés and B. Recht. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9:717–772, 2009.
- Kai-Yang Chiang, Cho-Jui Hsieh, and Inderjit S Dhillon. Matrix completion with noisy side information. In *Advances in Neural Information Processing Systems*, volume 28, page 3447–3455, 2015.
- D. M. Christopher, R. Prabhakar, and S. Hinrich. *Introduction to Information Retrieval*. Cambridge University Press, New York, 2008.
- J. M. Hernández-Lobato, N. Houlsby, and Z. Ghahramani. Probabilistic matrix factorization with non-random missing data. In *International Conference on Machine Learning*, pages 1512–1520. PMLR, 2014.
- D. Hsu, S. M. Kakade, and T. Zhang. Robust matrix decomposition with sparse corruptions. *IEEE Transactions on Information Theory*, 57(11):7221–7234, 2011.

- Y. Hu, Y. Koren, and C. Volinsky. Collaborative filtering for implicit feedback datasets. In *IEEE International Conference on Data Mining*, pages 263–272. IEEE, 2008.
- O. Klopp. Noisy low-rankmatrix completion with general sampling distribution. *Bernoulli*, 20:282–303, 2014.
- V. Koltchinskii, K. Lounici, and A. B. Tsybakov. Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *The Annals of Statistics*, 39:2302–2329, 2011.
- Michel Ledoux. *The Concentration of Measure Phenomenon*, volume 89. American Mathematical Soc., 2001.
- D. Liang, L. Charlin, J. McInerney, and D. M. Blei. Modeling user exposure in recommendation. In *International Conference on World Wide Web*, pages 951–961, 2016.
- Po-Ling Loh and Martin J Wainwright. High-dimensional regression with noisy and missing data: Provable guarantees with nonconvexity. *The Annals of Statistics*, 40(3):1637–1664, 2012.
- X. Mao, Wong R. K. W, and Chen S. X. Matrix completion under low-rank missing mechanism. *ArXiv*, abs/1812.07813, 2018.
- X. Mao, Chen S. X., and Wong R. K. W. Matrix completion with covariate information. *Journal of the American Statistical Association*, 114:198–210, 2019.
- B. M. Marlin and R. S. Zemel. Collaborative prediction and ranking with non-random missing data. In *ACM Conference on Recommender Systems*, pages 5–12, 2009.
- R. Mazumder, T. Hastie, and R. Tibshirani. Spectral regularization algorithms for learning large incomplete matrices. *Journal of Machine Learning Research*, 11:2287–2322, 2010.
- W. Miao and E. J. Tchetgen Tchetgen. On varieties of doubly robust estimators under missingness not at random with a shadow variable. *Biometrika*, 103(2):475–482, 2016.
- S. Negahban and M. J. Wainwright. Restricted strong convexity and weighted matrix completion: Optimal bounds with noise. *Journal of Machine Learning Research*, 13:1665–1697, 2012.
- R. Pan, Y. Zhou, B. Cao, N. N. Liu, R. Lukose, M. Scholz, and Q. Yang. One-class collaborative filtering. In *International Conference on Data Mining*, pages 502–511. IEEE, 2008.
- Gilles Pisier. *The Volume of Convex Bodies and Banach Space Geometry*, volume 94. Cambridge University Press, 1999.
- B. Recht. A simpler approach to matrix completion. *Journal of Machine Learning Research*, 12:3413–3430, 2011.
- G. Robin, H. T. Wai, J. Josse, O. Klopp, and É. Moulines. Low-rank interaction with sparse additive effects model for large data frames. In *Advances in Neural Information Processing Systems*, pages 5501–5511, 2018.

- G. Robin, O. Klopp, J. Josse, É. Moulines, and R. Tibshirani. Main effects and interactions in mixed and incomplete data frames. *Journal of the American Statistical Association*, 115(531):1292–1303, 2020.
- A. Rohde and A. B. Tsybakov. Estimation of high-dimensional low-rank matrices. *The Annals of Statistics*, 39:887–930, 2011.
- Nathan Srebro and Russ R Salakhutdinov. Collaborative filtering in a non-uniform world: Learning with the weighted trace norm. In *Advances in Neural Information Processing Systems*, volume 23, page 2056–2064, 2010.
- H. Steck. Training and testing of recommender systems on data missing not at random. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 713–722, 2010.
- S. A. Van De Geer, P. Bühlmann, et al. On the conditions used to prove oracle results for the lasso. *Electronic Journal of Statistics*, 3:1360–1392, 2009.
- A. Van Der Vaart and J. A. Wellner. *Weak Convergence and Empirical Processes*. Springer, 2000.
- Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019. doi: 10.1017/9781108627771.
- Miao Xu, Rong Jin, and Zhi-Hua Zhou. Speedup matrix completion with side information: Application to multi-label learning. In *Advances in Neural Information Processing Systems*, volume 26, page 2301–2309, 2013.
- J. Zhao and J. Shao. Semiparametric pseudo-likelihoods in generalized linear models with nonignorable missing data. *Journal of the American Statistical Association*, 110(512):1577–1590, 2015.
- Y. Zhu, X. Shen, and C. Ye. Personalized prediction and sparsity pursuit in latent factor models. *Journal of the American Statistical Association*, 111:241–252, 2016.