

# Doubly Robust Inference in Causal Latent Factor Models

Alberto Abadie      Anish Agarwal  
MIT                      Columbia

Raaz Dwivedi      Abhin Shah  
Cornell Tech      MIT

February 18, 2024

## Abstract

This article introduces a new framework for estimating average treatment effects under unobserved confounding in modern data-rich environments featuring large numbers of units and outcomes. The proposed estimator is doubly robust, combining outcome imputation, inverse probability weighting, and a novel cross-fitting procedure for matrix completion. We derive finite-sample and asymptotic guarantees, and show that the error of the new estimator converges to a mean-zero Gaussian distribution at a parametric rate. Simulation results demonstrate the practical relevance of the formal properties of the estimators analyzed in this article.

## 1. Introduction

This article presents a novel framework for the estimation of average treatment effects in modern data-rich environments in the presence of unobserved confounding. Modern data-rich environments are characterized by repeated measurements of outcomes, such as clinical metrics or purchase history, across a substantial number of units—be it patients in medical contexts or customers in online retail. As an example, consider an internet-retail platform where customers interact with various product categories. For each consumer-category pair, the platform makes decisions to either offer a discount or not, and records whether the consumer purchased a product in the category. Given an observational dataset capturing such interactions, our objective is to infer the causal effect of offering the discount on consumer purchase behavior. More specifically, we aim to infer two kinds of treatment effects: (a) tailored to product categories, the average impact of the discount on a product across consumers, and (b) tailored to consumers, the average impact of the discount on a consumer across product categories. This task is challenging due to unobserved confounding that may cause spurious associations between discount allocation and product purchase.

---

Abadie: [abadie@mit.edu](mailto:abadie@mit.edu). Agarwal: [aa5194@columbia.edu](mailto:aa5194@columbia.edu). Dwivedi: [dwivedi@cornell.edu](mailto:dwivedi@cornell.edu). Shah: [abhin@mit.edu](mailto:abhin@mit.edu).

There are two widely used approaches for treatment effect estimation: outcome-based methods and assignment-based methods. Outcome-based methods operate by imputing the missing potential outcomes for each consumer-product category pair. This process involves predicting whether a consumer, who received a discount, would have made the purchase without the discount (i.e., the potential outcome without discount), and conversely, if a consumer who did not receive the discount would have purchased the product had they received the discount (i.e., the potential outcome with discount). Assignment-based methods predict the probability with which a consumer is offered the discount on a product category, and inversely weight the observed outcomes by these estimated probabilities.

A substantial and influential body of literature has explored outcome-based methods, particularly in settings where all confounding factors are measured (see, e.g., Cochran, 1968; Rosenbaum and Rubin, 1983; Angrist, 1998; Abadie and Imbens, 2006, among many others). Imputing potential outcomes in the presence of unobserved confounders poses a more complex challenge, and the existing literature devoted to this problem is relatively small. In this context, a commonly adopted framework is the latent factor framework (Bai and Ng, 2002; Bai, 2009), wherein each element of the large-dimensional outcome vector is influenced by the same low-dimensional vector of unobserved confounders. A closely related approach is the technique of matrix completion (see, e.g., Chatterjee, 2015; Athey et al., 2021; Bai and Ng, 2021; Agarwal et al., 2023a; Dwivedi et al., 2022a) which has found widespread applications in recommendation systems and panel data models.

In this article, we propose a doubly-robust estimator (see Bang and Robins, 2005; Chernozhukov et al., 2018) of average treatment effects in the presence of unobserved confounding. This estimator leverages information on both the outcome process and the treatment assignment mechanism under a latent factor framework. It combines outcome imputation and inverse probability weighting with a new cross-fitting approach for matrix completion. We show that the proposed doubly-robust estimator has better finite-sample guarantees than alternative outcome-based and assignment-based estimators. Furthermore, the doubly-robust estimator is approximately Gaussian, asymptotically unbiased, and converges at a parametric rate, under provably valid error rates for matrix completion, irrespective of other properties of the matrix completion algorithm used for estimation, making it relatively agnostic to the specific matrix completion used.

**Terminology and notation.** For any real number  $b \in \mathbb{R}$ ,  $\lfloor b \rfloor$  is the greatest integer less than or equal to  $b$ . For any positive integer  $b$ ,  $[b]$  denotes the set of integers from 1 to  $b$ , i.e.,  $[b] \triangleq \{1, \dots, b\}$ . We use  $c$  to denote any generic universal constant, whose value may change between instances. For any  $c > 0$ ,  $m(c) = \max\{c, \sqrt{c}\}$  and  $\ell_c = \log(2/c)$ . For any two deterministic sequences  $a_n$  and  $b_n$  where  $b_n$  is positive,  $a_n = O(b_n)$  means that there exist a finite  $c > 0$  and a finite  $n_0 > 0$  such that  $|a_n| \leq c b_n$  for all  $n \geq n_0$ . Similarly,  $a_n = o(b_n)$  means that for every  $c > 0$ , there exists a finite  $n_0 > 0$  such that  $|a_n| < c b_n$  for all  $n \geq n_0$ . For a sequence of random variables  $x_n$  and a sequence of positive constants  $b_n$ ,  $x_n = O_p(b_n)$  means that the sequence  $|x_n/b_n|$  is stochastically bounded, i.e., for every  $\epsilon > 0$ , there exists a finite  $\delta > 0$  and a finite  $n_0 > 0$  such that  $\mathbb{P}(|x_n/b_n| > \delta) < \epsilon$  for all  $n \geq n_0$ . Similarly,  $x_n = o_p(b_n)$  means that the sequence  $|x_n/b_n|$  converges to zero in probability, i.e., for every

$\epsilon > 0$  and  $\delta > 0$ , there exists a finite  $n_0 > 0$  such that  $\mathbb{P}(|x_n/b_n| > \delta) < \epsilon$  for all  $n \geq n_0$ .

A mean-zero random variable  $x$  is subGaussian if there exists some  $b > 0$  such that  $\mathbb{E}[\exp(sx)] \leq \exp(b^2 s^2/2)$  for all  $s \in \mathbb{R}$ . Then, the subGaussian norm of  $x$  is given by  $\|x\|_{\psi_2} = \inf\{t > 0 : \mathbb{E}[\exp(x^2/t^2)] \leq 2\}$ . A mean-zero random variable  $x$  is subExponential if there exist some  $b_1, b_2 > 0$  such that  $\mathbb{E}[\exp(sx)] \leq \exp(b_1^2 s^2/2)$  for all  $-1/b_2 < s < 1/b_2$ . Then, the subExponential norm of  $x$  is given by  $\|x\|_{\psi_1} = \inf\{t > 0 : \mathbb{E}[\exp(|x|/t)] \leq 2\}$ . Let  $\text{Uniform}(a, b)$  denote the uniform distribution over the interval  $[a, b]$  for  $a, b \in \mathbb{R}$  such that  $a < b$ . Let  $\mathcal{N}(\mu, \sigma^2)$  denote the Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ .

For a vector  $u \in \mathbb{R}^n$ , we denote its  $t^{\text{th}}$  coordinate by  $u_t$  and its 2-norm  $\|u\|_2$ . For a matrix  $U \in \mathbb{R}^{n_1 \times n_2}$ , we denote the element in  $i^{\text{th}}$  row and  $j^{\text{th}}$  column by  $u_{i,j}$ , the  $i^{\text{th}}$  row by  $U_{i,\cdot}$ , the  $j^{\text{th}}$  column by  $U_{\cdot,j}$ , the largest eigenvalue by  $\lambda_{\max}(U)$ , and the smallest by  $\lambda_{\min}(U)$ . Given a set of indices  $\mathcal{R} \subseteq [n_1]$  and  $\mathcal{C} \subseteq [n_2]$ ,  $U_{\mathcal{I}} \in \mathbb{R}^{|\mathcal{R}| \times |\mathcal{C}|}$  is a sub-matrix of  $U$  corresponding to the entries in  $\mathcal{I} \triangleq \mathcal{R} \times \mathcal{C}$ . Further, we denote the Frobenius norm by  $\|U\|_{\text{F}} \triangleq (\sum_{i \in [n_1], j \in [n_2]} u_{i,j}^2)^{1/2}$ , the  $L_{1,2}$  norm by  $\|U\|_{1,2} \triangleq \max_{j \in [n_2]} (\sum_{i \in [n_1]} u_{i,j}^2)^{1/2}$ , the  $L_{2,\infty}$  norm by  $\|U\|_{2,\infty} \triangleq \max_{i \in [n_1]} (\sum_{j \in [n_2]} u_{i,j}^2)^{1/2}$ , and the maximum norm by  $\|U\|_{\max} \triangleq \max_{i \in [n_1], j \in [n_2]} |u_{i,j}|$ . Given two matrices  $U, V \in \mathbb{R}^{n_1 \times n_2}$ , the operators  $\odot$  and  $\oslash$  denote element-wise multiplication and division, respectively, i.e.,  $t_{i,j} = u_{i,j} \cdot v_{i,j}$  when  $T = U \odot V$ , and  $t_{i,j} = u_{i,j}/v_{i,j}$  when  $T = U \oslash V$ . When  $V$  is a binary matrix, i.e.,  $V \in \{0, 1\}^{n_1 \times n_2}$ , the operator  $\otimes$  is defined such that  $t_{i,j} = u_{i,j}$  if  $v_{i,j} = 1$  and  $t_{i,j} = ?$  if  $v_{i,j} = 0$  for  $T = U \otimes V$ . Given two matrices  $U \in \mathbb{R}^{n_1 \times n_2}$  and  $V \in \mathbb{R}^{n_1 \times n_3}$ , the operator  $*$  denotes the Khatri-Rao product (or column-wise product) of  $U$  and  $V$ , i.e.,  $T = U * V \in \mathbb{R}^{n_1 \times n_2 n_3}$  such that  $t_{i,j} = u_{i,j-n_2\bar{j}} \cdot v_{i,1+\bar{j}}$  where  $\bar{j} = \lfloor (j-1)/n_2 \rfloor$ . For random objects  $U$  and  $V$ ,  $U \perp\!\!\!\perp V$  means that  $U$  is independent of  $V$ .

## 2. Setup

Consider a setting with  $N$  units and  $M$  measurements per unit. For each unit-measurement pair  $i \in [N]$  and  $j \in [M]$ , we observe a treatment assignment  $a_{i,j} \in \{0, 1\}$  and the value of the outcome  $y_{i,j} \in \mathbb{R}$  under the treatment assignment. For the ease of exposition, we focus on binary treatments. However, our framework can be easily generalized to multi-ary treatments.

We operate within the Neyman-Rubin potential outcomes framework and denote the potential outcome for unit  $i \in [N]$  and measurement  $j \in [M]$  under treatment  $a \in \{0, 1\}$  by  $y_{i,j}^{(a)} \in \mathbb{R}$ . Here, it is implicitly assumed that the potential outcome for any unit  $i$  and measurement  $j$  does not depend on the treatment assignment for any other unit-measurement pair, i.e., there are no spillover effects across units or measurements. In the context of online retail data, the assumption of no spillovers across measurements is justified if the cross-elasticity of demand across product categories,  $j$ , is low. The observed outcomes depend on the potential outcomes and the treatment assignments,

$$y_{i,j} = y_{i,j}^{(0)}(1 - a_{i,j}) + y_{i,j}^{(1)}a_{i,j}, \quad (1)$$

for all  $i \in [N]$  and  $j \in [M]$ .

## 2.1. Sources of stochastic variation

In the setup of this article, each unit  $j \in [N]$  is characterized by a set of unknown parameters,  $\{(\theta_{i,j}^{(0)}, \theta_{i,j}^{(1)}, p_{i,j}) \in \mathbb{R}^2 \times [0, 1]\}_{j \in [M]}$ , which we treat as fixed. Potential outcomes and treatment assignments are generated as follows: for all  $i \in [N], j \in [M]$ , and  $a \in \{0, 1\}$ ,

$$y_{i,j}^{(a)} = \theta_{i,j}^{(a)} + \varepsilon_{i,j}^{(a)} \quad (2)$$

and

$$a_{i,j} = p_{i,j} + \eta_{i,j}, \quad (3)$$

where  $\varepsilon_{i,j}^{(a)}$  and  $\eta_{i,j}$  are mean-zero random variables, and

$$\eta_{i,j} = \begin{cases} -p_{i,j} & \text{with probability } 1 - p_{i,j} \\ 1 - p_{i,j} & \text{with probability } p_{i,j}. \end{cases} \quad (4)$$

It follows that  $\theta_{i,j}^{(a)}$  is the mean of the potential outcome  $y_{i,j}^{(a)}$ , and  $p_{i,j}$  is the unknown assignment probability or *latent propensity score*. The matrices  $\Theta^{(0)} \triangleq \{\theta_{i,j}^{(0)}\}_{i \in [N], j \in [M]}$ ,  $\Theta^{(1)} \triangleq \{\theta_{i,j}^{(1)}\}_{i \in [N], j \in [M]}$ , and  $P \triangleq \{p_{i,j}\}_{i \in [N], j \in [M]}$  collect all mean potential outcomes and assignment probabilities. Then, the matrices  $E^{(0)} \triangleq \{\varepsilon_{i,j}^{(0)}\}_{i \in [N], j \in [M]}$ ,  $E^{(1)} \triangleq \{\varepsilon_{i,j}^{(1)}\}_{i \in [N], j \in [M]}$ , and  $W \triangleq \{\eta_{i,j}\}_{i \in [N], j \in [M]}$  capture all sources of randomness in potential outcomes and treatment assignments.

Our setup allows  $\Theta^{(0)}, \Theta^{(1)}$  to be *arbitrarily* associated with  $P$ , inducing unobserved confounding. The identification restrictions made in Section 4 imply that  $\Theta^{(0)}, \Theta^{(1)}$ , and  $P$  include all confounding factors, and require  $(\varepsilon_{i,j}^{(0)}, \varepsilon_{i,j}^{(1)}) \perp\!\!\!\perp \eta_{i,j}$ .

## 2.2. Target causal estimand

For any given measurement  $j \in [M]$ , we aim to estimate the effect of the treatment averaged over all units,

$$\text{ATE}_{\cdot,j} \triangleq \mu_{\cdot,j}^{(1)} - \mu_{\cdot,j}^{(0)} \quad (5)$$

where

$$\mu_{\cdot,j}^{(a)} \triangleq \frac{1}{N} \sum_{i \in [N]} \theta_{i,j}^{(a)}.$$

It is straightforward to adapt the methods in this article to the estimation of alternative parameters, like the average treatment effect across measurements for each unit  $i$ , or the estimation of treatment effects over a subset of the units,  $S \subset [N]$ .

## 3. Estimation

In this section, we propose an estimator that uses the treatment assignment matrix  $A$  and the observed outcomes matrix  $Y$  to estimate the target causal estimand  $\{\text{ATE}_{\cdot,j}\}_{j \in [M]}$ , where

$$Y \triangleq \{y_{i,j}\}_{i \in [N], j \in [M]} \quad \text{and} \quad A \triangleq \{a_{i,j}\}_{i \in [N], j \in [M]}.$$

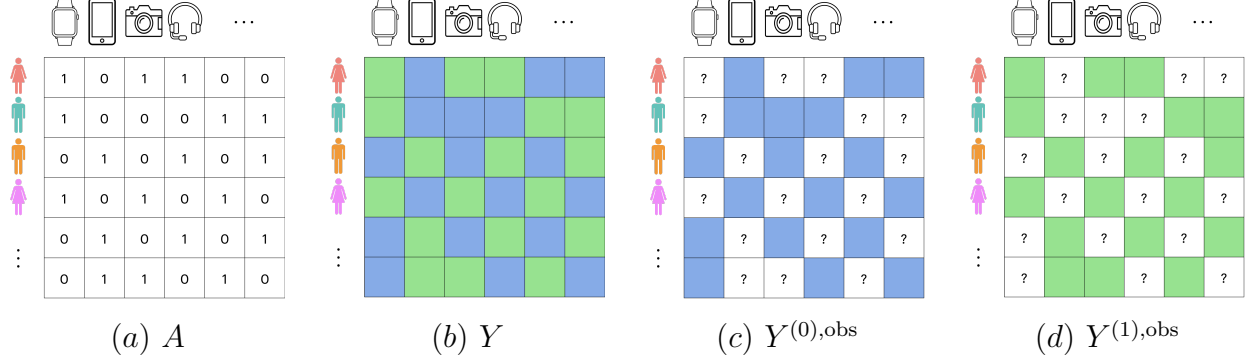


Figure 1: Schematic of the treatment assignment matrix  $A$ , the observed outcomes matrix  $Y$  (where green and blue fills indicate observations under  $a = 1$  and  $a = 0$ , respectively), and the observed component of the potential outcomes matrices, i.e.,  $Y^{(0),\text{obs}}$  and  $Y^{(1),\text{obs}}$  (where ? indicates a missing value). All matrices are  $N \times M$  where  $N$  is the number of customers and  $M$  is the number of products.

Our estimator leverages matrix completion as a key subroutine. We start with a brief overview of matrix completion below.

### 3.1. Matrix completion: A primer

Consider a matrix of parameters  $T \in \mathbb{R}^{N \times M}$ . While  $T$  is unobserved, we observe the matrix  $S \in \{\mathbb{R}, ?\}^{N \times M}$  where ? denotes a missing value. The relationship between  $S$  and  $T$  is given by

$$S = (T + H) \otimes F, \quad (6)$$

where  $H \in \mathbb{R}^{N \times M}$  represents a matrix of noise,  $F \in \{0, 1\}^{N \times M}$  is a masking matrix, and the operator  $\otimes$  is as defined in Section 1. A matrix completion algorithm, denoted by  $\text{MC}$ , takes the matrix  $S$  as its input, and returns an estimate for the matrix  $T$ , which we denote by  $\hat{T}$  or  $\text{MC}(S)$ . In other words,  $\text{MC}$  produces an estimate of a matrix from noisy observations of a subset of all the elements of the matrix.

The matrix completion literature is rich with algorithms  $\text{MC}$  that provide error guarantees, namely bounds on  $\|\text{MC}(S) - T\|$  for a suitably chosen norm/metric  $\|\cdot\|$ , under a variety of assumptions on the triplet  $(T, H, F)$ . Typical assumptions are (i)  $T$  is low-rank, (ii) the entries of  $H$  are independent, mean-zero and sub-Gaussian random variables, and (iii) the entries of  $F$  are independent Bernoulli random variables. Though matrix completion is commonly associated with the imputation of missing values, a typically underappreciated aspect is that it also denoises the observed matrix. Even when each entry of  $S$  is observed,  $\text{MC}(S)$  subtracts the effects of  $H$  from  $S$ , i.e., it performs matrix denoising. Refer to Nguyen et al. (2019) for a survey of various matrix completion algorithms.

### 3.2. Key building blocks

We now define and express matrices that are related to the quantities of interest  $\Theta^{(0)}$ ,  $\Theta^{(1)}$ , and  $P$  in a form similar to Eq. (6). See Figure 1 for a visual depiction of these matrices.

- **Outcomes:** Let  $Y^{(0),\text{obs}} = Y \otimes (\mathbf{1} - A) \in \{\mathbb{R}, ?\}^{N \times M}$  be a matrix with  $(i, j)$ -th entry equal to  $y_{i,j}$  if  $a_{i,j} = 0$  and equal to  $?$ , otherwise. Here,  $\mathbf{1}$  is the  $N \times M$  matrix with all entries equal to one. Analogously, let  $Y^{(1),\text{obs}} = Y \otimes A \in \{\mathbb{R}, ?\}^{N \times M}$  be a matrix with  $(i, j)$ -th entry equal to  $y_{i,j}$  if  $a_{i,j} = 1$  and equal to  $?$ , otherwise. In other words,  $Y^{(0),\text{obs}}$  and  $Y^{(1),\text{obs}}$  capture the observed components of  $\{y_{i,j}^{(0)}\}_{i \in [N], j \in [M]}$  and  $\{y_{i,j}^{(1)}\}_{i \in [N], j \in [M]}$ , respectively, with missing entries denoted by  $?$ . Then, we can write

$$Y^{(0),\text{obs}} = (\Theta^{(0)} + E^{(0)}) \otimes (\mathbf{1} - A) \quad \text{and} \quad Y^{(1),\text{obs}} = (\Theta^{(1)} + E^{(1)}) \otimes A. \quad (7)$$

- **Treatments:** From Eq. (3), we can write

$$A = (P + W),$$

as all the entries in  $A$  are observed. Building on the earlier discussion, the application of matrix completion yields the following estimates:

$$\hat{\Theta}^{(0)} = \text{MC}(Y^{(0),\text{obs}}), \quad \hat{\Theta}^{(1)} = \text{MC}(Y^{(1),\text{obs}}), \quad \text{and} \quad \hat{P} = \text{MC}(A), \quad (8)$$

where the algorithm  $\text{MC}$  may vary for  $\hat{\Theta}^{(0)}$ ,  $\hat{\Theta}^{(1)}$ , and  $\hat{P}$ . Because all entries of  $A$  are observed,  $\text{MC}(A)$  denoises  $A$  but does not need to impute missing entries. From Eq. (7) and Eq. (8), it follows that  $\hat{\Theta}^{(0)}$  and  $\hat{\Theta}^{(1)}$  depend on  $A$  and  $Y$ , whereas  $\hat{P}$  depends only on  $A$ .

In this section, we deliberately leave the matrix completion algorithm  $\text{MC}$  as a “black-box”. In Section 4, we establish finite-sample and asymptotic guarantees for our proposed estimator, contingent on specific properties for  $\text{MC}$ . In Section 5, we propose a novel end-to-end matrix completion algorithm that satisfies these properties.

Given matrix completion estimates of  $(\hat{\Theta}^{(0)}, \hat{\Theta}^{(1)}, \hat{P})$ , we formulate two preliminary estimators for  $\text{ATE}_{\cdot,j}$ : (i) an outcome imputation estimator, which uses  $\hat{\Theta}^{(0)}$  and  $\hat{\Theta}^{(1)}$  only, and (ii) an inverse probability weighting estimator, which uses  $\hat{P}$  only. Then, we combine these to obtain a doubly robust estimator of  $\text{ATE}_{\cdot,j}$ .

**Outcome imputation (OI) estimator.** Let  $\hat{\theta}_{i,j}^{(a)}$  denote the  $(i, j)$ -th entry of  $\hat{\Theta}^{(a)}$  for  $i \in [N]$ ,  $j \in [M]$ , and  $a \in \{0, 1\}$ . The OI estimator for  $\text{ATE}_{\cdot,j}$  is defined as follows:

$$\widehat{\text{ATE}}_{\cdot,j}^{\text{OI}} \triangleq \hat{\mu}_{\cdot,j}^{(1,\text{OI})} - \hat{\mu}_{\cdot,j}^{(0,\text{OI})}, \quad (9)$$

where

$$\hat{\mu}_{\cdot,j}^{(a,\text{OI})} \triangleq \frac{1}{N} \sum_{i \in [N]} \hat{\theta}_{i,j}^{(a)} \quad \text{for } a \in \{0, 1\}.$$

That is, the OI estimator is obtained by taking the difference of the average value of the  $j$ -th column of the estimates  $\hat{\Theta}^{(0)}$  and  $\hat{\Theta}^{(1)}$ . The quality of the OI estimator depends on how well  $\hat{\Theta}^{(0)}$  and  $\hat{\Theta}^{(1)}$  approximate the mean potential outcome matrices  $\Theta^{(0)}$  and  $\Theta^{(1)}$ , respectively.

**Inverse probability weighting (IPW) estimator.** Let  $\hat{p}_{i,j}$  denote the  $(i, j)$ -th entry of  $\hat{P}$  for  $i \in [N]$  and  $j \in [M]$ . The IPW estimate for  $\text{ATE}_{\cdot,j}$  is defined as follows:

$$\widehat{\text{ATE}}_{\cdot,j}^{\text{IPW}} \triangleq \hat{\mu}_{\cdot,j}^{(1,\text{IPW})} - \hat{\mu}_{\cdot,j}^{(0,\text{IPW})}, \quad (10)$$

where

$$\hat{\mu}_{\cdot,j}^{(0,\text{IPW})} \triangleq \frac{1}{N} \sum_{i \in [N]} \frac{y_{i,j}(1 - a_{i,j})}{1 - \hat{p}_{i,j}} \quad \text{and} \quad \hat{\mu}_{\cdot,j}^{(1,\text{IPW})} \triangleq \frac{1}{N} \sum_{i \in [N]} \frac{y_{i,j}a_{i,j}}{\hat{p}_{i,j}}.$$

That is, the IPW estimator is obtained by taking the difference of the average value of the  $j$ -th column of the matrices  $Y^{(0),\text{obs}}$  and  $Y^{(1),\text{obs}}$ , replacing unobserved entries with zeros, and weighting each outcome by the inverse of the estimated assignment probability to account for confounding. The quality of the IPW estimate depends on how well  $\hat{P}$  approximates the probability matrix  $P$ .

The matrix completion-based OI and IPW estimators in Eq. (9) and Eq. (10) have the same form as the classical OI and IPW estimators, which are derived for settings where all confounders are observed (e.g., Imbens and Rubin, 2015). In contrast to the classical setting, our framework is one with unmeasured confounding.

### 3.3. Doubly robust (DR) estimator

The DR estimate for  $\text{ATE}_{\cdot,j}$  combines the estimates  $\hat{\Theta}^{(0)}$ ,  $\hat{\Theta}^{(1)}$ , and  $\hat{P}$  from Eq. (8). It is defined as follows:

$$\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}} \triangleq \hat{\mu}_{\cdot,j}^{(1,\text{DR})} - \hat{\mu}_{\cdot,j}^{(0,\text{DR})}, \quad (11)$$

where

$$\hat{\mu}_{\cdot,j}^{(0,\text{DR})} \triangleq \frac{1}{N} \sum_{i \in [N]} \hat{\theta}_{i,j}^{(0,\text{DR})} \quad \text{with} \quad \hat{\theta}_{i,j}^{(0,\text{DR})} \triangleq \hat{\theta}_{i,j}^{(0)} + (y_{i,j} - \hat{\theta}_{i,j}^{(0)}) \frac{1 - a_{i,j}}{1 - \hat{p}_{i,j}},$$

and

$$\hat{\mu}_{\cdot,j}^{(1,\text{DR})} \triangleq \frac{1}{N} \sum_{i \in [N]} \hat{\theta}_{i,j}^{(1,\text{DR})} \quad \text{with} \quad \hat{\theta}_{i,j}^{(1,\text{DR})} \triangleq \hat{\theta}_{i,j}^{(1)} + (y_{i,j} - \hat{\theta}_{i,j}^{(1)}) \frac{a_{i,j}}{\hat{p}_{i,j}}. \quad (12)$$

In Section 4, we prove that  $\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}}$  consistently estimates  $\text{ATE}_{\cdot,j}$  as long as either  $(\hat{\Theta}^{(0)}, \hat{\Theta}^{(1)})$  is consistent for  $(\Theta^{(0)}, \Theta^{(1)})$  or  $\hat{P}$  is consistent for  $P$ , i.e., it is doubly robust. Furthermore, we show that the DR estimator provides superior finite sample guarantees than the OI and IPW estimators, and that it satisfies a central limit theorem at a parametric rate under weak conditions on the convergence rate of the matrix completion routine. Using simulated data, Figure 2 demonstrates the improved performance of DR, relative to OI and IPW. Despite substantial biases observed in both OI and IPW estimates, the error of the DR estimate demonstrates a mean-zero Gaussian distribution. We provide a detailed description of the simulation setup in Section 6.

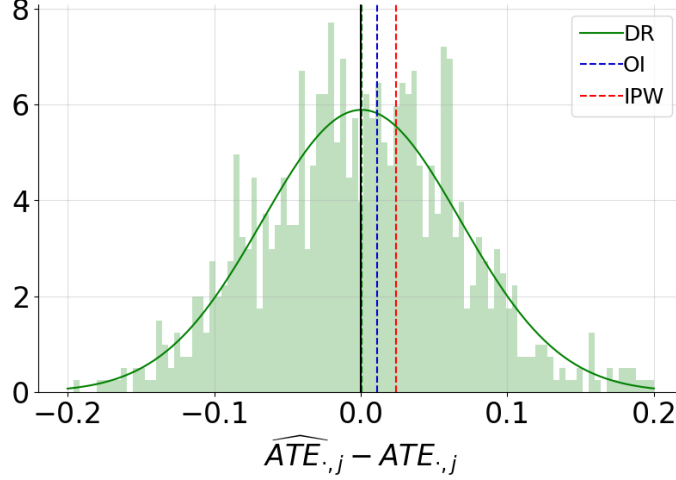


Figure 2: Empirical illustration of the convergence of the error of the doubly robust (DR) estimator to a mean-zero Gaussian distribution. The histogram represents  $\widehat{ATE}_{\cdot,j}^{DR} - ATE_{\cdot,j}$  and the curve represents the (best) fitted Gaussian distribution. Histogram counts are normalized so that the area under the histogram integrates to one. Unlike DR, the outcome imputation (OI) and inverse probability weighting (IPW) estimators have non-trivial biases, as evidenced by the means of the distributions in dashed green, blue, and red, respectively. We provide details of the simulations, including the data-generating process, in Section 6.

## 4. Main Results

This section presents the formal results of the article. Section 4.1 details assumptions, Section 4.2 discusses finite-sample guarantees, and Section 4.3 presents a central limit theorem for  $\widehat{ATE}_{\cdot,j}^{DR}$ .

### 4.1. Assumptions

**Requirements on data generating process.** We make two assumptions on how the data is generated. First, we impose a positivity condition on the assignment probabilities.

**Assumption 1** (Positivity). *The unknown assignment probability matrix  $P$  is such that*

$$\lambda \leq p_{i,j} \leq 1 - \lambda, \quad (13)$$

*for all  $i \in [N]$  and  $j \in [M]$ , where  $0 < \lambda \leq 1/2$  is a constant.*

Assumption 1 requires that the propensity score for each unit-outcome pair is bounded away from 0 and 1, implying that any unit-item pair can be assigned either of the two treatments. An analogous assumption is pervasive in causal inference models that assume observed confounding. For simplicity of exposition and to avoid notational clutter, Assumption 1 requires Eq. (13) for all outcomes,  $j \in [M]$ . However, it is only necessary that Eq. (13) holds for the outcomes of interest,  $j$ , for which  $ATE_{\cdot,j}$  is estimated. Our framework leverages the



availability of a large number of outcomes to control for the confounding effect of latent variables. In practical applications, however,  $\text{ATE}_{\cdot,j}$  may be estimated for a select group of those outcomes. For example, in synthetic control settings (Abadie et al., 2010),  $\text{ATE}_{\cdot,j}$  is estimated only for post-treatment outcomes. In that case, the positivity assumption applies only for the selected subset of outcomes for which  $\text{ATE}_{\cdot,j}$  is estimated.

Next, we formalize the requirements on the noise variables.

**Assumption 2** (Zero-mean, independent, and subGaussian noise).

- (a)  $(W, E^{(0)}, E^{(1)})$  have zero mean entries,
- (b)  $W \perp\!\!\!\perp (E^{(0)}, E^{(1)})$ ,
- (c) All entries of  $W$  are mutually independent,
- (d)  $\{(\varepsilon_{i,j}^{(0)}, \varepsilon_{i,j}^{(1)}) : i \in [N]\}$  are mutually independent (across  $i$ ) for every  $j \in [M]$ , and
- (e) Each entry of  $E^{(0)}$  and  $E^{(1)}$  has subGaussian norm bounded by a constant  $\bar{\sigma}$ .

Assumption 2(a) defines  $(\Theta^{(0)}, \Theta^{(1)}, P)$  as the means of the potential outcomes and treatment assignment in Eqs. (2) and (3). Assumption 2(b) implies that  $(\Theta^{(0)}, \Theta^{(1)}, P)$  capture all confounding factors. Assumption 2(c) imposes independence across units and measurements in the noise  $W$ . Assumption 2(d) imposes independence across units in the noise  $(E^{(0)}, E^{(1)})$ , for every measurement. Finally, Assumption 2(e) is mild and useful to derive finite-sample guarantees. For the central limit theorem in Section 4.3, subGaussianity could be disposed of by restricting the moments of  $\varepsilon_{i,j}^{(0)}$  and  $\varepsilon_{i,j}^{(1)}$ . Note that Assumption 2 does not restrict the dependence between  $\varepsilon_{i,j}^{(0)}$  and  $\varepsilon_{i,j}^{(1)}$ .

**Requirements on matrix completion estimators.** First, we assume the estimate  $\hat{P}$  is consistent with Assumption 1.

**Assumption 3.** The estimated probability matrix  $\hat{P}$  is such that

$$\bar{\lambda} \leq \hat{p}_{i,j} \leq 1 - \bar{\lambda},$$

for all  $i \in [N]$  and  $j \in [M]$ , where  $0 < \bar{\lambda} \leq \lambda$ .

Assumption 3 is achieved by truncating entries of  $\hat{P}$  to the range  $[\bar{\lambda}, 1 - \bar{\lambda}]$ . Second, our theoretical analysis requires independence between certain sub-matrices of the estimates  $(\hat{P}, \hat{\Theta}^{(0)}, \hat{\Theta}^{(1)})$  from Eq. (8), and the noise matrices  $(W, E^{(0)}, E^{(1)})$ . We formally state this independence condition as an assumption below.

**Assumption 4.** There exists partitions  $(\mathcal{R}_0, \mathcal{R}_1)$  of the units in  $[N]$  and  $(\mathcal{C}_0, \mathcal{C}_1)$  of the measurements  $[M]$ , such that each unit  $i \in [N]$  is assigned to  $\mathcal{R}_0$  or  $\mathcal{R}_1$  with equal probability, each measurement  $j \in [M]$  is assigned to  $\mathcal{C}_0$  or  $\mathcal{C}_1$  with equal probability, and for each block  $\mathcal{I} \in \mathcal{P} \triangleq \{\mathcal{R}_i \times \mathcal{C}_j : i, j \in \{0, 1\}\}$ ,

$$\hat{P}_{\mathcal{I}}, \hat{\Theta}_{\mathcal{I}}^{(0)}, \hat{\Theta}_{\mathcal{I}}^{(1)} \perp\!\!\!\perp W_{\mathcal{I}} \tag{14}$$

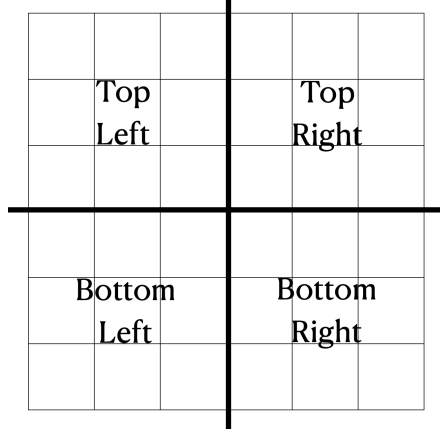


Figure 3: A matrix  $S$  partitioned into four blocks when  $\mathcal{R}_0 = \lfloor N/2 \rfloor$  and  $\mathcal{C}_0 = \lfloor M/2 \rfloor$  in Assumption 4, i.e.,  $\mathcal{P} = \{\text{Top Left}, \text{Top Right}, \text{Bottom Left}, \text{Bottom Right}\}$ .

and

$$\hat{P}_{\mathcal{I}} \perp\!\!\!\perp E_{\mathcal{I}}^{(0)}, E_{\mathcal{I}}^{(1)}. \quad (15)$$

Without loss of generality, suppose  $\mathcal{R}_0 = \lfloor \lfloor N/2 \rfloor \rfloor$  and  $\mathcal{C}_0 = \lfloor \lfloor M/2 \rfloor \rfloor$ . Figure 3 provides a schematic of the corresponding block partition  $\mathcal{P}$ . Eq. (14) requires that within each of the four blocks in  $\mathcal{P}$ , mean potential outcomes estimators and the assignment probability estimators are independent of the sub-matrix of  $W$  for the same block. Assumption 2(b) implies that Eq. (15) holds provided  $\hat{P}$  is a function of  $A$  only, as is the case for the matrix completion procedure in Eq. (8). Analogous conditions appear in the literature on doubly robust estimation under observed confounding (e.g., Definition 3.1 in Chernozhukov et al., 2018). Specifically, in that context, Chernozhukov et al. (2018) split the available data into  $K$ -folds, and require estimates of propensities and outcomes in each fold to be independent of the noise in that fold. Section 5 provides a way to ensure Assumption 4 holds for any MC algorithm using a cross-fitting procedure as long as Assumption 2 holds.

**Matrix completion error rates.** The formal guarantees in this section depend on the normalized  $L_{1,2}$  norms of the errors in estimating the unknown parameters  $(\Theta^{(0)}, \Theta^{(1)}, P)$ . We use the following notation for these errors:

$$\mathcal{E}(\hat{P}) \triangleq \frac{\|\hat{P} - P\|_{1,2}}{\sqrt{N}} \quad \text{and} \quad \mathcal{E}(\hat{\Theta}) \triangleq \sum_{a \in \{0,1\}} \mathcal{E}(\hat{\Theta}^{(a)}), \quad (16)$$

where

$$\mathcal{E}(\hat{\Theta}^{(a)}) = \frac{\|\hat{\Theta}^{(a)} - \Theta^{(a)}\|_{1,2}}{\sqrt{N}}.$$

A variety of matrix completion algorithms deliver  $\mathcal{E}(\hat{P}) = O_p(\min\{N, M\}^{-\alpha})$  and  $\mathcal{E}(\hat{\Theta}) = O_p(\min\{N, M\}^{-\beta})$ , where  $0 < \alpha, \beta \leq 1/2$ . Throughout, our notation primarily tracks

dependence on  $N$ . We say that these normalized errors achieve the parametric rate when they have the same rate as  $O_p(N^{-1/2})$ . Section 5 explicitly characterizes  $\alpha$  and  $\beta$  under low-rank assumptions on  $(\Theta^{(0)}, \Theta^{(1)})$  and  $P$  for a particular matrix completion algorithm.

## 4.2. Non-asymptotic guarantees

The first main result of this section provides both a non-asymptotic error bound and an asymptotic consistency result for  $\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}} - \text{ATE}_{\cdot,j}$  in terms of the errors  $\mathcal{E}(\hat{P})$  and  $\mathcal{E}(\hat{\Theta})$  in Eq. (16).

**Theorem 1 (Finite Sample Guarantees for DR).** *Suppose Assumptions 1 to 4 hold. Fix  $\delta \in (0, 1)$  and  $j \in [M]$ . Then, with probability at least  $1 - \delta$ , we have*

$$|\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}} - \text{ATE}_{\cdot,j}| \leq \text{Err}_N^{\text{DR}}, \quad (17)$$

where

$$\text{Err}_N^{\text{DR}} \triangleq \frac{2}{\lambda} \left[ \mathcal{E}(\hat{\Theta}) \cdot \mathcal{E}(\hat{P}) + \left( \frac{\sqrt{c\ell_{\delta/12}}}{\sqrt{\ell_1}} \mathcal{E}(\hat{\Theta}) + 2\bar{\sigma} \sqrt{c\ell_{\delta/12}} + \frac{2\bar{\sigma}m(c\ell_{\delta/12})}{\sqrt{\ell_1}} \right) \cdot \frac{1}{\sqrt{N}} \right], \quad (18)$$

for  $m(c)$  and  $\ell_c$  as defined in Section 1. Therefore, as  $N \rightarrow \infty$ , if either (i)  $\mathcal{E}(\hat{P}) = o_p(1)$ ,  $\mathcal{E}(\hat{\Theta}) = O_p(1)$ , or (ii)  $\mathcal{E}(\hat{\Theta}) = o_p(1)$ ,  $\mathcal{E}(\hat{P}) = O_p(1)$ , it holds that

$$\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}} \xrightarrow{p} \text{ATE}_{\cdot,j}, \quad (19)$$

for all  $j \in [M]$ .

The proof of Theorem 1 is given in Appendix B. Eqs. (17) and (18) bound the absolute error of the DR estimator by the rate of  $\mathcal{E}(\hat{\Theta})(\mathcal{E}(\hat{P}) + N^{-0.5}) + N^{-0.5}$ . When  $\mathcal{E}(\hat{P})$  is lower bounded at the parametric rate of  $N^{-0.5}$ ,  $\text{Err}_N^{\text{DR}}$  has the same rate as  $\mathcal{E}(\hat{P})\mathcal{E}(\hat{\Theta}) + N^{-0.5}$ .

**Doubly robust behavior of  $\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}}$ .** The error rate of  $\mathcal{E}(\hat{P})\mathcal{E}(\hat{\Theta}) + N^{-0.5}$  immediately reveals that the DR estimate is doubly robust with respect to the error in estimating the mean potential outcomes  $(\Theta^{(0)}, \Theta^{(1)})$  and the assignment probabilities  $P$ . First, the error  $\text{Err}_N^{\text{DR}}$  decays at a parametric rate of  $O_p(N^{-0.5})$  as long as the product of error rates,  $\mathcal{E}(\hat{P})\mathcal{E}(\hat{\Theta})$ , decays as  $O_p(N^{-0.5})$ . As a result,  $\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}}$  can exhibit a parametric error rate even when neither the mean potential outcomes nor the assignment probabilities are estimated at a parametric rate. Second,  $\text{Err}_N^{\text{DR}}$  decays to zero as long as either of  $\mathcal{E}(\hat{P})$  or  $\mathcal{E}(\hat{\Theta})$  decays to 0. Hence,  $\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}}$  is consistent as long as either the mean potential outcomes or the assignment probabilities are estimated consistently.

We next compare the performance of DR estimator with the OI and IPW estimators from Eqs. (9) and (10), respectively. Towards this goal, we characterize the  $\text{ATE}_{\cdot,j}$  estimation error of  $\widehat{\text{ATE}}_{\cdot,j}^{\text{OI}}$  in terms of  $\mathcal{E}(\hat{\Theta})$  and of  $\widehat{\text{ATE}}_{\cdot,j}^{\text{IPW}}$  in terms of  $\mathcal{E}(\hat{P})$ .

**Proposition 1 (Finite Sample Guarantees for OI and IPW).** *Fix any  $j \in [M]$ . For OI, we have*

$$|\widehat{\text{ATE}}_{\cdot,j}^{\text{OI}} - \text{ATE}_{\cdot,j}| \leq \text{Err}_N^{\text{OI}} \triangleq \mathcal{E}(\widehat{\Theta}). \quad (20)$$

*For IPW, suppose Assumptions 1 to 4 hold. Define  $\|\Theta\|_{\max} \triangleq \sum_{a \in \{0,1\}} \|\Theta^{(a)}\|_{\max}$ , and fix any  $\delta \in (0, 1)$ . Then, with probability at least  $1 - \delta$ , we have*

$$|\widehat{\text{ATE}}_{\cdot,j}^{\text{IPW}} - \text{ATE}_{\cdot,j}| \leq \text{Err}_N^{\text{IPW}}, \quad (21)$$

where

$$\text{Err}_N^{\text{IPW}} \triangleq \frac{2}{\lambda} \left[ \|\Theta\|_{\max} \cdot \mathcal{E}(\widehat{P}) + \left( \frac{\sqrt{c\ell_{\delta/12}}}{\sqrt{\ell_1}} \|\Theta\|_{\max} + 2\bar{\sigma} \sqrt{c\ell_{\delta/12}} + \frac{2\bar{\sigma}m(c\ell_{\delta/12})}{\sqrt{\ell_1}} \right) \cdot \frac{1}{\sqrt{N}} \right],$$

for  $m(c)$  and  $\ell_c$  as defined in Section 1.

The proofs of Eq. (20) and Eq. (21) are given in Appendices D and E, respectively. Proposition 1 implies that in an asymptotic sequence with bounded  $\|\Theta\|_{\max}$ , OI and IPW attain the parametric rate  $O_p(N^{-0.5})$  provided  $\mathcal{E}(\widehat{\Theta})$  and  $\mathcal{E}(\widehat{P})$  are  $O_p(N^{-0.5})$ , respectively. The next corollary compares these error rates with those obtained for the DR estimator in Theorem 1.

**Corollary 1 (Gains of DR over OI and IPW).** *Suppose Assumptions 1 to 4 hold. Consider an asymptotic sequence such that  $\|\Theta\|_{\max}$  is bounded. If  $\mathcal{E}(\widehat{P}) = O_p(N^{-\alpha})$  and  $\mathcal{E}(\widehat{\Theta}) = O_p(N^{-\beta})$  for  $0 \leq \alpha, \beta \leq 0.5$ , then*

$$\text{Err}_N^{\text{OI}} = O_p(N^{-\beta}), \quad \text{Err}_N^{\text{IPW}} = O_p(N^{-\alpha}),$$

and

$$\text{Err}_N^{\text{DR}} = O_p(N^{-\min\{\alpha+\beta, 0.5\}}).$$

Corollary 1 demonstrates that the DR estimate's error decay rate is consistently superior to that of the OI and IPW estimates across a variety of regimes for  $\alpha, \beta$ . Specifically, the error  $\text{Err}_N^{\text{DR}}$  scales strictly faster than both  $\text{Err}_N^{\text{OI}}$  and  $\text{Err}_N^{\text{IPW}}$  if the estimation errors of  $\widehat{\Theta}^{(0)}$ ,  $\widehat{\Theta}^{(1)}$ , and  $\widehat{P}$  converge slower than at the parametric rate  $O_p(N^{-1/2})$ . When the estimation errors of  $\widehat{\Theta}^{(0)}$ ,  $\widehat{\Theta}^{(1)}$ , and  $\widehat{P}$  all decay at a parametric rate, OI, IPW, and DR estimation errors decay also at a parametric rate.

### 4.3. Gaussian approximation

The next theorem, proven in Appendix C, establishes a Gaussian approximation for  $\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}}$  under mild conditions on error rates  $\mathcal{E}(\widehat{P})$  and  $\mathcal{E}(\widehat{\Theta})$ .

**Theorem 2 (Asymptotic Normality for DR).** *Suppose Assumptions 1 to 4 and the following conditions hold,*

(C1)  $\mathcal{E}(\hat{P}) = O_p(s_N)$  and  $\mathcal{E}(\hat{\Theta}) = O_p(t_N)$  where the sequences  $s_N$  and  $t_N$  are  $o(1)$ .

(C2)  $\mathcal{E}(\hat{P})\mathcal{E}(\hat{\Theta}) = o_p(N^{-1/2})$ .

(C3) Let  $\sigma_{i,j}^{(0)}$  and  $\sigma_{i,j}^{(1)}$  be the standard deviations of  $\varepsilon_{i,j}^{(0)}$  and  $\varepsilon_{i,j}^{(1)}$ , respectively. The sequence

$$\bar{\sigma}_j^2 \triangleq \frac{1}{N} \sum_{i \in [N]} \frac{(\sigma_{i,j}^{(1)})^2}{p_{i,j}} + \frac{1}{N} \sum_{i \in [N]} \frac{(\sigma_{i,j}^{(0)})^2}{1 - p_{i,j}}, \quad (22)$$

is bounded away from zero as  $N$  increases.

Then, for all  $j \in [M]$ ,

$$\sqrt{N}(\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}} - \text{ATE}_{\cdot,j})/\bar{\sigma}_j \xrightarrow{d} \mathcal{N}(0, 1), \quad (23)$$

as  $N \rightarrow \infty$ .

Theorem 2 describes two simple requirements on the estimated  $\hat{P}$  and  $(\hat{\Theta}^{(0)}, \hat{\Theta}^{(1)})$ , under which  $\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}}$  exhibits an asymptotic Gaussian distribution centered at  $\text{ATE}_{\cdot,j}$ . Condition (C1) requires that the estimation errors of  $\hat{P}$  and  $(\hat{\Theta}^{(0)}, \hat{\Theta}^{(1)})$  converge to zero in probability. Condition (C2) requires that the product of the errors decays sufficiently fast, at a rate  $o_p(N^{-1/2})$ , ensuring that the bias of the normalized estimator in Eq. (23) converges to zero. Condition (C2) is similar to conditions in the literature on doubly-robust estimation of average treatment effects under observed confounding (e.g., Assumption 5.1 in Chernozhukov et al., 2018). Specifically, in that context, Chernozhukov et al. (2018) assume that the product of propensity estimation error and outcome regression error decays faster than  $N^{-1/2}$ .

**Black-box asymptotic normality.** We emphasize Theorem 2 applies to any matrix completion algorithm MC as long as conditions (C1) and (C2) are satisfied. This property arises because the bias is dominated by the product of  $\mathcal{E}(\hat{P})$  and  $\mathcal{E}(\hat{\Theta})$ , which can be shown to be  $o_p(N^{-1/2})$  for a broad class of MC algorithms under mild assumptions on  $(P, \Theta^{(0)}, \Theta^{(1)})$ . On the other hand, achieving such black-box asymptotic normality results for OI or IPW estimates is challenging, as their bias scales with individual error rates  $\mathcal{E}(\hat{\Theta})$  and  $\mathcal{E}(\hat{P})$ , respectively, which are typically lower bounded at the parametric rate of  $N^{-0.5}$ . Our simulations in Section 6 corroborate these theoretical findings.

## 5. Matrix Completion with Cross-Fitting

In this section, we introduce a novel algorithm designed to construct estimators  $(\hat{\Theta}^{(0)}, \hat{\Theta}^{(1)}, \hat{P})$  that adhere to Assumption 4 and satisfy conditions (C1) and (C2) in Theorem 2. We first explain why traditional matrix completion algorithms fail to deliver the properties required by Assumption 4. We then present **Cross-Fitted-MC**, a meta-algorithm that takes any matrix completion algorithm and uses it to construct  $(\hat{\Theta}^{(0)}, \hat{\Theta}^{(1)}, \hat{P})$  that satisfy Assumption 4. Finally, we describe **Cross-Fitted-SVD**, an end-to-end algorithm obtained by combining

**Cross-Fitted-MC** with the singular value decomposition (SVD)-based algorithm of Bai and Ng (2021), and establish that it also satisfies conditions (C1) and (C2) in Theorem 2.

**Traditional matrix completion.** Estimators  $(\hat{\Theta}^{(0)}, \hat{\Theta}^{(1)}, \hat{P})$  obtained from existing matrix completion algorithms need not satisfy Assumption 4. In particular, using the entire assignment matrix  $A$  to estimate each element of  $P$  typically results in a violation of  $\hat{P}_{\mathcal{I}} \perp\!\!\!\perp W_{\mathcal{I}}$  in Eq. (14), as each entry of  $\hat{P}$  is allowed to depend on the entire noise matrix  $W$ . For example, in spectral methods (e.g., Nguyen et al., 2019),  $\hat{P}$  is a function of the SVD of the entire matrix  $A$ , and

$$\hat{p}_{i,j} \not\perp\!\!\!\perp a_{i',j'}, \quad (24)$$

for all  $(i, j), (i', j') \in [N] \times [M]$  in general, which implies that for every  $\mathcal{I} \subseteq [N] \times [M]$ ,  $\hat{P}_{\mathcal{I}} \not\perp\!\!\!\perp W_{\mathcal{I}}$ . Similarly, in matching methods such as nearest neighbors (Li et al., 2019),  $\hat{P}$  is a function of the matches/neighbors estimated from the entire matrix  $A$ . Dependence structures such as  $\hat{p}_{i,j} \not\perp\!\!\!\perp a_{i,j}$  for any  $i, j \in [N] \times [M]$ —which is weaker than Eq. (24)—are enough to violate the  $\hat{P}_{\mathcal{I}} \perp\!\!\!\perp W_{\mathcal{I}}$  requirement in Eq. (14).

Likewise, the requirement  $\hat{\Theta}_{\mathcal{I}}^{(0)}, \hat{\Theta}_{\mathcal{I}}^{(1)} \perp\!\!\!\perp W_{\mathcal{I}}$  in Eq. (14) can be violated, because  $\hat{\Theta}^{(0)}$  and  $\hat{\Theta}^{(1)}$  depend respectively on  $Y^{(0),\text{obs}}$  and  $Y^{(1),\text{obs}}$ , which themselves depend on the entire matrix  $A$ .

### 5.1. Cross-Fitted-MC: A meta-cross-fitting algorithm for matrix completion

We now introduce **Cross-Fitted-MC**, a cross-fitting approach that modifies any MC algorithm to produce  $(\hat{\Theta}^{(0)}, \hat{\Theta}^{(1)}, \hat{P})$  that satisfy Assumption 4. Recall the setup from Section 3.1: Given an observation matrix  $S \in \{\mathbb{R}, ?\}^{N \times M}$ , a matrix completion algorithm MC produces an estimate  $\hat{T} = \text{MC}(S) \in \mathbb{R}^{N \times M}$  of a matrix of interest  $T$ , where  $S$  and  $T$  are related via Eq. (6). With this background, we now describe the **Cross-Fitted-MC** meta-algorithm.

1. The inputs are (i) a matrix completion algorithm MC, (ii) an observation matrix  $S \in \{\mathbb{R}, ?\}^{N \times M}$ , and (iii) a block partition  $\mathcal{P}$  of the set  $[N] \times [M]$  into four blocks as in Assumption 4.
2. For each block  $\mathcal{I} \in \mathcal{P}$ , construct  $\hat{T}_{\mathcal{I}}$  by applying MC on  $S \otimes \mathbf{1}^{-\mathcal{I}}$  where  $\mathbf{1}^{-\mathcal{I}} \in \mathbb{R}^{N \times M}$  denotes a masking matrix with  $(i, j)$ -th entry equal to 0 if  $(i, j) \in \mathcal{I}$  and 1 otherwise, and the operator  $\otimes$  is as defined in Section 1. In other words,

$$\hat{T}_{\mathcal{I}} = \bar{T}_{\mathcal{I}} \quad \text{where} \quad \bar{T} = \text{MC}(S \otimes \mathbf{1}^{-\mathcal{I}}).$$

3. Return  $\hat{T} \in \mathbb{R}^{N \times M}$  obtained by collecting together  $\{\hat{T}_{\mathcal{I}}\}_{\mathcal{I} \in \mathcal{P}}$ , with each entry in its original position.

We represent this meta-algorithm succinctly as below:

$$\hat{T} = \text{Cross-Fitted-MC}(\text{MC}, S, \mathcal{P}).$$

In summary, **Cross-Fitted-MC** produces an estimator  $\hat{T}$  such that for each block  $\mathcal{I} \in \mathcal{P}$ , the sub-matrix  $\hat{T}_{\mathcal{I}}$  is constructed only using the entries of  $S$  corresponding to the remaining three blocks of  $\mathcal{P}$ . See Figure 4 for a visualization of  $S \otimes \mathbf{1}^{-\mathcal{I}}$ . The following result, proven in Appendix F.1, establishes  $(\hat{\Theta}^{(0)}, \hat{\Theta}^{(1)}, \hat{P})$  generated by **Cross-Fitted-MC** satisfy Assumption 4.

**Proposition 2 (Guarantees for Cross-Fitted-MC).** *Suppose Assumption 2 holds. Let MC be any matrix completion algorithm and  $\mathcal{P}$  be any block partition of the set  $[N] \times [M]$  into four blocks as in Assumption 4. Let*

$$\hat{\Theta}^{(0)} = \text{Cross-Fitted-MC}(\text{MC}, Y^{(0),\text{obs}}, \mathcal{P}), \quad (25)$$

$$\hat{\Theta}^{(1)} = \text{Cross-Fitted-MC}(\text{MC}, Y^{(1),\text{obs}}, \mathcal{P}), \quad (26)$$

$$\hat{P} = \text{Cross-Fitted-MC}(\text{MC}, A, \mathcal{P}), \quad (27)$$

where  $Y^{(0),\text{obs}}$  and  $Y^{(1),\text{obs}}$  are defined in Eq. (7). Then, Assumption 4 holds.

A host of MC algorithms are designed to de-noise and impute missing entries of matrices under random patterns of missingness; the most common missingness pattern studied is where each entry has the same probability of being missing, independent of everything else. In contrast, **Cross-Fitted-MC** generates patterns where all entries in one block are deterministically missing, as in Figure 4. A recent strand of research on the interplay between matrix completion methods and causal inference models—specifically, within the synthetic controls framework—has contributed matrix completion algorithms that allow for block missingness (see, e.g., Athey et al., 2021; Agarwal et al., 2021; Bai and Ng, 2021; Agarwal et al., 2023b; Arkhangelsky et al., 2021; Agarwal et al., 2023a; Dwivedi et al., 2022a,b). However, it is a challenge to apply known theoretical guarantees for these methods to the setting in this article because of: (i) the use of cross-fitting—which creates blocks where all observations are missing—and (ii) outside of the completely-missing blocks, there can still be missing observations with heterogeneous probabilities of missingness. In the next section, we show how to modify any MC algorithm designed for block missingness patterns so that it can be applied to our setting with cross-fitting and heterogeneous probabilities of missingness outside the folds. For concreteness, we work with the Tall-Wide matrix completion algorithm of Bai and Ng (2021).

## 5.2. The Cross-Fitted-SVD algorithm

**Cross-Fitted-SVD** is an end-to-end MC algorithm obtained by instantiating the **Cross-Fitted-MC** meta-algorithm with the Tall-Wide algorithm of Bai and Ng (2021), which we denote as TW. For completeness, we detail the TW algorithm in Section 5.2.1, and then use it to describe **Cross-Fitted-SVD** in Section 5.2.2.

### 5.2.1. The TW algorithm of Bai and Ng (2021).

Bai and Ng (2021) propose TW to impute missing values in those matrices where there exists a set of rows and a set of columns without missing entries. More concretely, for

|  |        |  |   |       |   |
|--|--------|--|---|-------|---|
|  |        |  |   |       |   |
|  | Top    |  |   | Top   |   |
|  | Left   |  |   | Right |   |
|  |        |  |   |       |   |
|  |        |  | ? | ?     | ? |
|  | Bottom |  | ? | ?     | ? |
|  | Left   |  | ? | ?     | ? |
|  |        |  | ? | ?     | ? |

Figure 4: The matrix  $S \otimes \mathbf{1}^{\text{Bottom Right}}$  obtained from the matrix  $S$  in Figure 3 by masking the entries corresponding to the Bottom Right block with  $?$ .

any matrix  $S \in \{\mathbb{R}, ?\}^{N \times M}$ , let  $\mathcal{R}_{\text{obs}} \subseteq [N]$  and  $\mathcal{C}_{\text{obs}} \subseteq [M]$  denote the set of rows and columns, respectively, with all entries observed. Then, the block  $\mathcal{I} = \mathcal{R}_{\text{miss}} \times \mathcal{C}_{\text{miss}}$ , where  $\mathcal{R}_{\text{miss}} \triangleq [N] \setminus \mathcal{R}_{\text{obs}}$  and  $\mathcal{C}_{\text{miss}} \triangleq [M] \setminus \mathcal{C}_{\text{obs}}$ , is such that all the missing entries in  $S$  are a subset of it.

Given a rank hyper-parameter  $r \in [\min\{|\mathcal{R}_{\text{obs}}|, |\mathcal{C}_{\text{obs}}|\}]$ ,  $\text{TW}_r$  produces an estimate of  $T$  as follows:

1. Run SVD separately on  $S^{(\text{tall})} \triangleq S_{[N] \times \mathcal{C}_{\text{obs}}}$  and  $S^{(\text{wide})} \triangleq S_{\mathcal{R}_{\text{obs}} \times [M]}$ , i.e.,

$$\text{SVD}(S^{(\text{tall})}) = (U^{(\text{tall})} \in \mathbb{R}^{N \times \bar{r}_N}, \Sigma^{(\text{tall})} \in \mathbb{R}^{\bar{r}_N \times \bar{r}_N}, V^{(\text{tall})} \in \mathbb{R}^{|\mathcal{C}_{\text{obs}}| \times \bar{r}_N})$$

and

$$\text{SVD}(S^{(\text{wide})}) = (U^{(\text{wide})} \in \mathbb{R}^{|\mathcal{R}_{\text{obs}}| \times \bar{r}_M}, \Sigma^{(\text{wide})} \in \mathbb{R}^{\bar{r}_M \times \bar{r}_M}, V^{(\text{wide})} \in \mathbb{R}^{M \times \bar{r}_M})$$

where  $\bar{r}_N \triangleq \min\{N, |\mathcal{C}_{\text{obs}}|\}$  and  $\bar{r}_M \triangleq \min\{|\mathcal{R}_{\text{obs}}|, M\}$ . The columns of  $U^{(\text{tall})}$  and  $U^{(\text{wide})}$  are the left singular vectors of  $S^{(\text{tall})}$  and  $S^{(\text{wide})}$ , respectively, and the columns of  $V^{(\text{tall})}$  and  $V^{(\text{wide})}$  are the right singular vectors of  $S^{(\text{tall})}$  and  $S^{(\text{wide})}$ , respectively. The diagonal entries of  $\Sigma^{(\text{tall})}$  and  $\Sigma^{(\text{wide})}$  are the singular values of  $S^{(\text{tall})}$  and  $S^{(\text{wide})}$ , respectively, and the off-diagonal entries are zeros. This step of  $\text{TW}$  requires the existence of the fully observed blocks  $S^{(\text{tall})}$  and  $S^{(\text{wide})}$ , i.e.,  $\mathcal{R}_{\text{obs}}$  and  $\mathcal{C}_{\text{obs}}$  cannot be empty.

2. Let  $\tilde{V}^{(\text{tall})} \in \mathbb{R}^{|\mathcal{C}_{\text{obs}}| \times r}$  be the sub-matrix of  $V^{(\text{tall})}$  that keeps the columns corresponding to the  $r$  largest singular values only. Let  $\tilde{V}^{(\text{wide})} \in \mathbb{R}^{|\mathcal{C}_{\text{obs}}| \times r}$  be the sub-matrix of  $V^{(\text{wide})}$  that keeps the columns corresponding to the  $r$  largest singular values only and the rows corresponding to the indices in  $\mathcal{C}_{\text{obs}}$  only. Obtain a rotation matrix  $R \in \mathbb{R}^{r \times r}$  as follows:

$$R \triangleq \tilde{V}^{(\text{tall})\top} \tilde{V}^{(\text{wide})} (\tilde{V}^{(\text{wide})\top} \tilde{V}^{(\text{wide})})^{-1}.$$

That is,  $R$  is obtained by regressing  $\tilde{V}^{(\text{tall})}$  on  $\tilde{V}^{(\text{wide})}$ . In essence,  $R$  aligns the right singular vectors of  $S^{(\text{tall})}$  and  $S^{(\text{wide})}$  using the entries that are common between these two



matrices, i.e., the entries corresponding to indices  $\mathcal{R}_{\text{obs}} \times \mathcal{C}_{\text{obs}}$ . The formal guarantees of the TW algorithm remains unchanged if one alternatively regresses  $\tilde{V}^{(\text{wide})}$  on  $\tilde{V}^{(\text{tall})}$ , or uses the left singular vectors of  $S^{(\text{tall})}$  and  $S^{(\text{wide})}$  for alignment.

3. Let  $\bar{\Sigma}^{(\text{tall})} \in \mathbb{R}^{\bar{r}_N \times r}$  be the sub-matrix of  $\Sigma^{(\text{tall})}$  that keeps the columns corresponding to the  $r$  largest singular values only. Let  $\bar{V}^{(\text{wide})} \in \mathbb{R}^{M \times r}$  be the sub-matrix of  $V^{(\text{wide})}$  that keeps the columns corresponding to the  $r$  largest singular values only. Return  $\hat{T} \triangleq U^{(\text{tall})} \bar{\Sigma}^{(\text{tall})} R \bar{V}^{(\text{wide})\top}$  as an estimate for  $T$ .

### 5.2.2. *Cross-Fitted-SVD algorithm.*

1. The inputs are (i)  $A \in \mathbb{R}^{N \times M}$ , (ii)  $Y^{(a), \text{obs}} \in \{\mathbb{R}, ?\}^{N \times M}$  for  $a \in \{0, 1\}$ , and (iii) hyper-parameters  $r_1, r_2, r_3$ , and  $\bar{\lambda}$  such that  $r_1, r_2, r_3 \in [\min\{N, M\}]$  and  $0 < \bar{\lambda} \leq 1/2$ .
2. Choose a random partition  $(\mathcal{R}_0, \mathcal{R}_1)$  of  $[N]$  and  $(\mathcal{C}_0, \mathcal{C}_1)$  of  $[M]$  such that each  $i \in [N]$  is assigned to  $\mathcal{R}_0$  or  $\mathcal{R}_1$  with equal probability and each  $j \in [M]$  is assigned to  $\mathcal{C}_0$  or  $\mathcal{C}_1$  with equal probability. Construct the block partition  $\mathcal{P} \triangleq \{\mathcal{R}_i \times \mathcal{C}_j : i, j \in \{0, 1\}\}$ .
3. Return  $\hat{P} = \text{Proj}_{\bar{\lambda}}(\text{Cross-Fitted-MC}(\text{TW}_{r_1}, A, \mathcal{P}))$  where  $\text{Proj}_{\bar{\lambda}}(\cdot)$  projects each entry of its input to the interval  $[\bar{\lambda}, 1 - \bar{\lambda}]$ .
4. Define  $Y^{(0), \text{full}}$  as equal to  $Y^{(0), \text{obs}}$ , but with all missing entries in  $Y^{(0), \text{obs}}$  set to zero. Define  $Y^{(0), \text{obs}}$  analogously with respect to  $Y^{(1), \text{full}}$ .
5. Return  $\hat{\Theta}^{(0)} = \text{Cross-Fitted-MC}(\text{TW}_{r_2}, Y^{(0), \text{full}}, \mathcal{P}) \odot (\mathbf{1} - \hat{P})$ .
6. Return  $\hat{\Theta}^{(1)} = \text{Cross-Fitted-MC}(\text{TW}_{r_3}, Y^{(1), \text{full}}, \mathcal{P}) \odot \hat{P}$ .

We provide intuition on the key steps of the **Cross-Fitted-SVD** algorithm next.

**Computing  $\hat{P}$ .** The estimate  $\hat{P}$  comes from applying **Cross-Fitted-MC** with TW on  $A$  and truncating the entries of the resulting matrix to the range  $[\bar{\lambda}, 1 - \bar{\lambda}]$ , in accordance with Assumption 3. The TW sub-routine is directly applicable to  $A$ , because for any block  $\mathcal{I} = \mathcal{R}_i \times \mathcal{C}_j \in \mathcal{P}$  the masked matrix  $A \otimes \mathbf{1}^{-\mathcal{I}}$  has  $[N] \setminus \mathcal{R}_i$  fully observed rows and  $[M] \setminus \mathcal{C}_j$  fully observed columns. See Figure 5(a) for a visualization of  $A \otimes \mathbf{1}^{-\mathcal{I}}$ .

**Computing  $\hat{\Theta}^{(0)}$  and  $\hat{\Theta}^{(1)}$ .** The estimates  $\hat{\Theta}^{(0)}$  and  $\hat{\Theta}^{(1)}$  are constructed by applying **Cross-Fitted-MC** with TW on  $Y^{(0), \text{full}}$  and  $Y^{(1), \text{full}}$ , which do not have missing entries. TW is not directly applicable on  $Y^{(0), \text{obs}}$  and  $Y^{(1), \text{obs}}$ , as both matrices may not have any rows and columns that are fully observed. See Figure 5(b) and Figure 5(c) for visualizations of  $Y^{(0), \text{obs}} \otimes \mathbf{1}^{-\mathcal{I}}$  and  $Y^{(1), \text{obs}} \otimes \mathbf{1}^{-\mathcal{I}}$ , respectively. However, notice that

$$\mathbb{E}[Y^{(0), \text{full}}] = \mathbb{E}[Y \odot (\mathbf{1} - A)] = \Theta^{(0)} \odot (\mathbf{1} - P),$$

and

$$\mathbb{E}[Y^{(1), \text{full}}] = \mathbb{E}[Y \odot A] = \Theta^{(1)} \odot P.$$

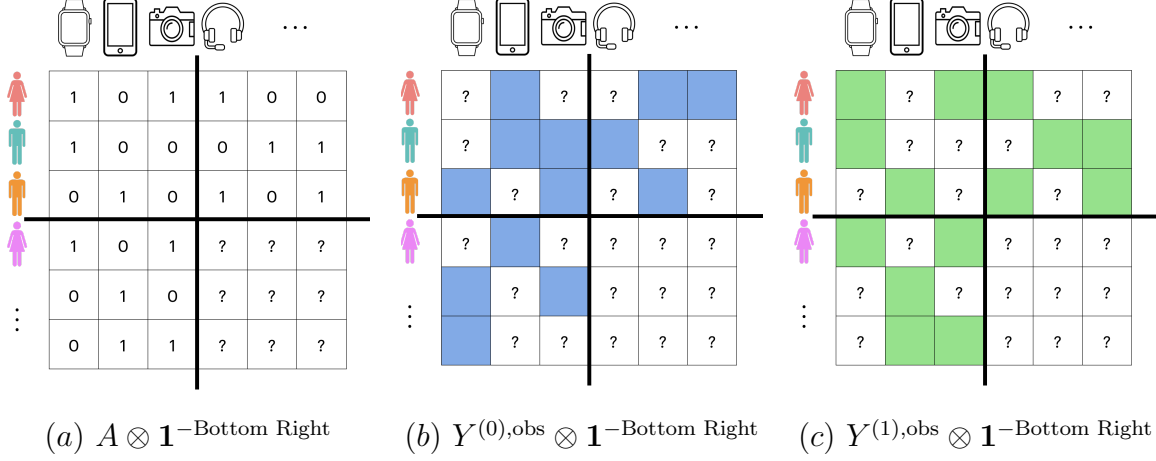


Figure 5: Panels (a), (b), and (c) illustrate the matrices  $A \otimes \mathbf{1}^{-\mathcal{I}}$ ,  $Y^{(0),\text{obs}} \otimes \mathbf{1}^{-\mathcal{I}}$ , and  $Y^{(1),\text{obs}} \otimes \mathbf{1}^{-\mathcal{I}}$  obtained from  $A$ ,  $Y^{(0),\text{obs}}$  and  $Y^{(1),\text{obs}}$ , respectively, for the block partition  $\mathcal{P}$  in Figure 3 and the block  $\mathcal{I} = \text{Bottom Right}$ . Unlike Panels (b) and (c), there exists rows and columns with all entries observed in Panel (a). To enable the application of TW for Panels (b) and (c), we replace missing entries in blocks Top Left, Top Right, and Bottom Left with zeros.

As a result,  $\text{MC}(Y^{(0),\text{full}})$  and  $\text{MC}(Y^{(1),\text{full}})$  provide estimates of  $\Theta^{(0)} \odot (\mathbf{1} - P)$  and  $\Theta^{(1)} \odot P$ , respectively—recall the discussion in Section 3.1. To estimate  $\Theta^{(0)}$  and  $\Theta^{(1)}$ , we divide the entries of  $\text{MC}(Y^{(0),\text{full}})$  and  $\text{MC}(Y^{(1),\text{full}})$  by the entries of  $(\mathbf{1} - \hat{P})$  and  $\hat{P}$ , respectively. Adjustments of this type for heterogeneous missingness probabilities have been previously explored in Ma and Chen (2019); Bhattacharya and Chatterjee (2022).

### 5.3. Theoretical guarantees for Cross-Fitted-SVD

To establish theoretical guarantees for **Cross-Fitted-SVD**, we adopt three assumptions from Bai and Ng (2021). The first assumption imposes a low-rank structure on the matrices  $P$ ,  $\Theta^{(0)}$ , and  $\Theta^{(1)}$ , namely that their entries are given by an inner product of latent factors.

**Assumption 5** (Linear latent factor model on the confounders). *There exist constants  $r_p, r_{\theta_0}, r_{\theta_1} \in [\min\{N, M\}]$  and a collection of latent factors*

$$U \in \mathbb{R}^{N \times r_p}, \quad V \in \mathbb{R}^{M \times r_p}, \quad U^{(a)} \in \mathbb{R}^{N \times r_{\theta_a}}, \quad \text{and} \quad V^{(a)} \in \mathbb{R}^{M \times r_{\theta_a}} \quad \text{for } a \in \{0, 1\},$$

*such that the unobserved confounders  $(\Theta^{(0)}, \Theta^{(1)}, P)$  satisfy the following factorization:*

$$P = UV^\top \quad \text{and} \quad \Theta^{(a)} = U^{(a)}V^{(a)\top} \quad \text{for } a \in \{0, 1\}. \quad (28)$$

Assumption 5 decomposes each of the unobserved confounders ( $P$ ,  $\Theta^{(0)}$ , and  $\Theta^{(1)}$ ) into low-dimensional unit-dependent latent factors ( $U$ ,  $U^{(0)}$ , and  $U^{(1)}$ ) and measurement-dependent latent factors ( $V$ ,  $V^{(0)}$ , and  $V^{(1)}$ ). In particular, every unit  $i \in [N]$  is associated with three low-dimensional factors: (i)  $U_{i,\cdot} \in \mathbb{R}^{r_p}$ , (ii)  $U_{i,\cdot}^{(0)} \in \mathbb{R}^{r_{\theta_0}}$ , and (iii)  $U_{i,\cdot}^{(1)} \in \mathbb{R}^{r_{\theta_1}}$ . Similarly, every

measurement  $j \in [M]$  is associated with three factors: (i)  $V_{i,\cdot} \in \mathbb{R}^{r_p}$ , (ii)  $V_{i,\cdot}^{(0)} \in \mathbb{R}^{r_{\theta_0}}$ , and (iii)  $V_{i,\cdot}^{(1)} \in \mathbb{R}^{r_{\theta_1}}$ . Such low-rank assumptions are standard in the matrix completion literature.

The second assumption requires that the factors that determine  $P$ ,  $\Theta^{(0)} \odot (\mathbf{1} - P)$ , and  $\Theta^{(1)} \odot P$  explain a sufficiently large amount of the variation in the data. This assumption is made on the factors of  $\Theta^{(0)} \odot (\mathbf{1} - P)$  and  $\Theta^{(1)} \odot P$  instead of  $\Theta^{(0)}$  and  $\Theta^{(1)}$  as the TW algorithm is applied on  $Y^{(0),\text{full}} = Y \odot (\mathbf{1} - A)$  and  $Y^{(1),\text{full}} = Y \odot A$ , instead of  $Y^{(0),\text{obs}}$  and  $Y^{(1),\text{obs}}$  (see steps 5 and 6 of **Cross-Fitted-SVD**). To determine the factors of  $\Theta^{(0)} \odot (\mathbf{1} - P)$  and  $\Theta^{(1)} \odot P$ , let

$$\bar{U} \triangleq [\mathbf{1}_N, -U] \in \mathbb{R}^{N \times (r_p+1)} \quad \text{and} \quad \bar{V} \triangleq [\mathbf{1}_M, -V] \in \mathbb{R}^{M \times (r_p+1)},$$

where  $\mathbf{1}_N \in \mathbb{R}^N$  and  $\mathbf{1}_M \in \mathbb{R}^M$  are vectors of all 1's. Then,

$$\Theta^{(0)} \odot (\mathbf{1} - P) = \bar{U}^{(0)} \bar{V}^{(0)\top} \quad \text{and} \quad \Theta^{(1)} \odot P = \bar{U}^{(1)} \bar{V}^{(1)\top}, \quad (29)$$

where  $\bar{U}^{(0)} \triangleq \bar{U} * U^{(0)} \in \mathbb{R}^{N \times r_{\theta_0}(r_p+1)}$ ,  $\bar{V}^{(0)} \triangleq \bar{V} * V^{(0)} \in \mathbb{R}^{M \times r_{\theta_0}(r_p+1)}$ ,  $\bar{U}^{(1)} \triangleq U * U^{(1)} \in \mathbb{R}^{N \times r_{\theta_1}r_p}$ , and  $\bar{V}^{(1)} \triangleq V * V^{(1)} \in \mathbb{R}^{M \times r_{\theta_1}r_p}$ , with the operator  $*$  denoting the row-wise Khatri-Rao product (see Section 1). We provide details of the derivation of these factors in Appendix F.2.3.

**Assumption 6** (Strong factors). *There exists a positive constant  $c$  such that*

$$\|U\|_{2,\infty} \leq c, \quad \|V\|_{2,\infty} \leq c, \quad \|U^{(a)}\|_{2,\infty} \leq c, \quad \text{and} \quad \|V^{(a)}\|_{2,\infty} \leq c \quad \text{for } a \in \{0, 1\}.$$

*Further, the matrices defined below are positive definite:*

$$\Sigma^U \triangleq \frac{U^\top U}{N}, \quad \Sigma^V \triangleq \frac{V^\top V}{M}, \quad \Sigma^{\bar{U}^{(a)}} \triangleq \frac{\bar{U}^{(a)\top} \bar{U}^{(a)}}{N}, \quad \text{and} \quad \Sigma^{\bar{V}^{(a)}} \triangleq \frac{\bar{V}^{(a)\top} \bar{V}^{(a)}}{M} \quad \text{for } a \in \{0, 1\}.$$

Assumption 6, a classic assumption in the literature on latent factor models, ensures that the factor structure is strong. Specifically, it ensures that each eigenvector of  $P$ ,  $\Theta^{(0)} \odot (\mathbf{1} - P)$ , and  $\Theta^{(1)} \odot P$  carries sufficiently large signal.

The subsequent assumption introduces additional conditions on the noise variables in Bai and Ng (2021) than those specified in Assumption 2.

**Assumption 7** (Weak dependence across measurements and independence across units).

(a)  $\sum_{j' \in [M]} |\mathbb{E}[\varepsilon_{i,j}^{(a)} \varepsilon_{i,j'}^{(a)}]| \leq c$  for every  $i \in [N]$ ,  $j \in [M]$ , and  $a \in \{0, 1\}$ , and

(b)  $\{E_{i,\cdot}^{(a)} : i \in [N]\}$  are mutually independent (across  $i$ ) for  $a \in \{0, 1\}$ .

For every  $a \in \{0, 1\}$ , Assumption 7(a) requires the noise  $E^{(a)}$  to exhibit only weak dependency across measurements and Assumption 7(b) requires the noise  $E^{(a)}$  to be independent across units. We are now ready to provide guarantees on the estimates produced by **Cross-Fitted-SVD**. The proof can be found in Appendix F.2.

**Proposition 3 (Guarantees for Cross-Fitted-SVD).** *Suppose Assumptions 1, 2, and 5 to 7 hold. Consider an asymptotic sequence such that  $\|\Theta\|_{\max}$  is bounded as both  $N$  and  $M$  increase. Let  $\hat{P}$ ,  $\hat{\Theta}^{(0)}$ , and  $\hat{\Theta}^{(1)}$  be the estimates returned by **Cross-Fitted-SVD** with  $r_1 = r_p$ ,  $r_2 = r_{\theta_0}(r_p + 1)$ ,  $r_3 = r_{\theta_1}r_p$ , and any  $\bar{\lambda}$  such that  $0 < \bar{\lambda} \leq \lambda$  with  $\lambda$  denoting the constant from Assumption 1. Then, as  $N, M \rightarrow \infty$ ,*

$$\mathcal{E}(\hat{P}) = O_p\left(\frac{1}{\sqrt{N}} + \frac{1}{\sqrt{M}}\right) \quad \text{and} \quad \mathcal{E}(\hat{\Theta}) = O_p\left(\frac{1}{\sqrt{N}} + \frac{1}{\sqrt{M}}\right).$$

Proposition 3 implies that the conditions (C1) and (C2) in Theorem 2 hold whenever  $N^{1/2}/M = o(1)$ . Then, the DR estimator from Eq. (11) constructed using the estimates  $\hat{\Theta}^{(0)}$ ,  $\hat{\Theta}^{(1)}$ , and  $\hat{P}$  returned by **Cross-Fitted-SVD** exhibits an asymptotic Gaussian distribution centered at the target causal estimand. Further, Proposition 3 implies that the estimation errors  $\mathcal{E}(\hat{P})$  and  $\mathcal{E}(\hat{\Theta})$  achieve the parametric rate whenever  $N/M = o(1)$ .

## 6. Simulations

This section reports simulation results on the performance of the DR estimator of Eq. (11) and the OI and IPW estimators of Eqs. (9) and (10), respectively. For convenience, we let  $N = M$ .

**Data Generating Process (DGP).** We now briefly describe the DGP for our simulations; details can be found in Appendix G. To generate,  $P$ ,  $\Theta^{(0)}$ , and  $\Theta^{(1)}$ , we use the latent factor model given in Eq. (28). To introduce unobserved confounding, we set the unit-specific latent factors to be the same across  $P$ ,  $\Theta^{(0)}$ , and  $\Theta^{(1)}$ , i.e.,  $U = U^{(0)} = U^{(1)}$ . The entries of  $U$  and the measurement-specific latent factors,  $V, V^{(0)}, V^{(1)}$  are each sampled independently from a uniform distribution. Further, the entries of the noise matrices  $E^{(0)}$  and  $E^{(1)}$  are sampled independently from a normal distribution, and the entries of  $W$  are sampled independently as per Eq. (4). Then,  $y_{i,j}^{(a)}$ ,  $a_{i,j}$ , and  $y_{i,j}$  are determined from Eqs. (1) to (3), respectively. The simulation generates  $P$ ,  $\Theta^{(0)}$ , and  $\Theta^{(1)}$  once. Then, given the fixed values of  $P$ ,  $\Theta^{(0)}$ , and  $\Theta^{(1)}$ , the simulation generates  $Q$  realizations of  $(Y, A)$ —that is, only the noise matrices  $E^{(0)}, E^{(1)}, W$  are resampled for each of the  $Q$  realizations. For each of these  $Q$  instances of the simulation,  $\hat{P}$ ,  $\hat{\Theta}^{(0)}$ , and  $\hat{\Theta}^{(1)}$  are obtained by applying the **Cross-Fitted-SVD** algorithm to the corresponding  $A$  and  $Y$  with the choice of hyper-parameters as in Proposition 3 and  $\bar{\lambda} = \lambda = 0.05$ . For each of the instances of the simulation we compute  $\text{ATE}_{\cdot,j}$  from Eq. (5), and  $\widehat{\text{ATE}}_{\cdot,j}^{\text{OI}}$ ,  $\widehat{\text{ATE}}_{\cdot,j}^{\text{IPW}}$  and  $\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}}$  from Eqs. (9) to (11). We set  $Q = 1000$ . While Proposition 3 assumes the bounds  $r_p$  and  $r_{\theta}$  on the ranks of the latent factors from Assumption 5 are constants, we relax this restriction in the simulations and allow  $r_p$  and  $r_{\theta}$  to scale with  $N$  in the simulations, as we note below.

**Results.** Figure 6 reports simulation results for  $N = 1000$ , with  $r_p = \lfloor N^{1/5} \rfloor$ ,  $r_{\theta} = \lfloor N^{1/4} \rfloor$  in Panel (a), and  $r_p = \lfloor N^{1/4} \rfloor$ ,  $r_{\theta} = \lfloor N^{1/5} \rfloor$  in Panel (b). Figure 2 in Section 3 reports simulation results for  $r_p = r_{\theta} = \lfloor N^{1/5} \rfloor$ . In each case, the figure shows a histogram of the distribution of  $\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}} - \text{ATE}_{\cdot,j}$  across simulation instances, along with the best fitting Gaussian distribution

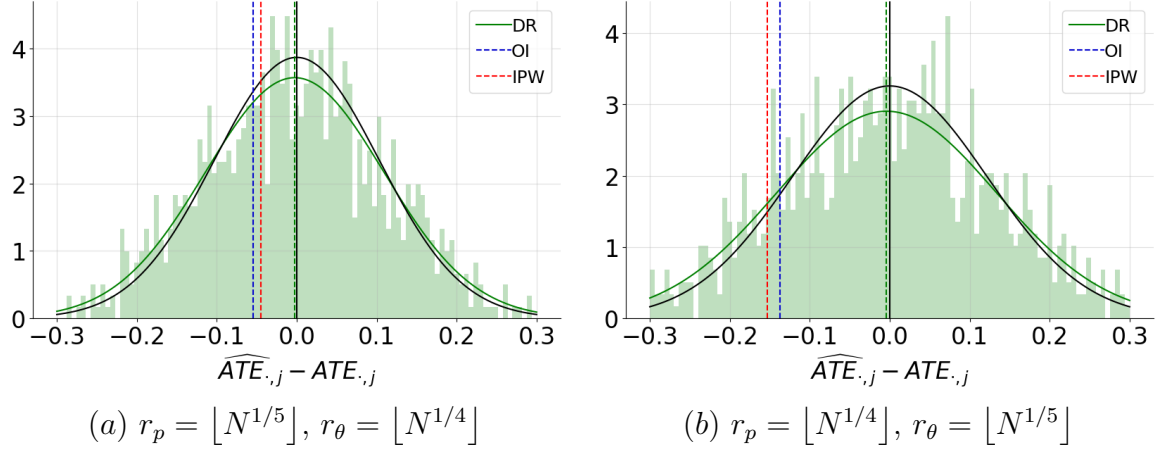


Figure 6: Empirical illustration of the asymptotic performance of DR as in Theorem 2. The histogram corresponds to the errors of 1000 independent instances of DR estimates, the green curve represents the (best) fitted Gaussian distribution, and the black curve represents the Gaussian approximation from Theorem 2. The dashed green, blue, and red lines represent the biases of DR, OI, and IPW estimators.

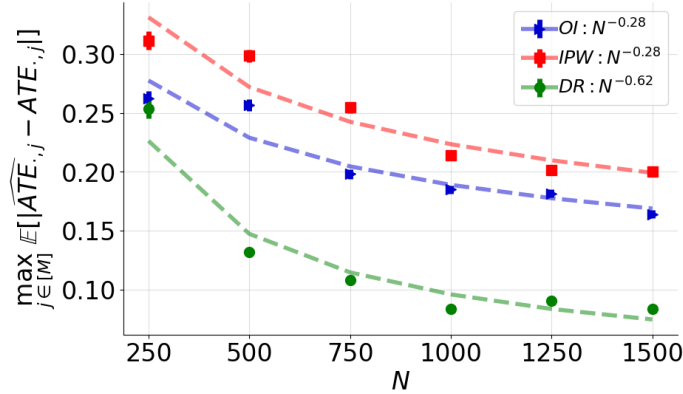


Figure 7: Comparison of OI, IPW, and DR in terms of finite sample performance as in Proposition 1. The estimates  $\widehat{ATE}_{.j}^{OI}$ ,  $\widehat{ATE}_{.j}^{IPW}$ , and  $\widehat{ATE}_{.j}^{DR}$  are obtained by taking an average over 1000 independent instances.

(green curve). The histogram counts are normalized so that the area under the histogram integrates to one. Figure 6 plots the Gaussian distribution in the result of Theorem 2 (black curve). The dashed blue, red and green lines in Figures 2 and 6 indicate the values of the means of the OI, IPW, and DR error, respectively, across simulation instances. For reference, we place a black solid line at zero and the black curve represents the Gaussian approximation from Theorem 2. The DR estimator has minimal bias and a close-to-Gaussian distribution. The biases of OI and IPW are non-negligible.

To further illustrate the different bias performance of the three estimators, Figure 7 reports

the maximum over  $j \in [M]$  of their respective mean absolute error estimates. For each  $j$ , the estimate of the mean absolute error of OI, IPW, and DR is the average of  $|\widehat{\text{ATE}}_{\cdot,j}^{\text{OI}} - \text{ATE}_{\cdot,j}|$ ,  $|\widehat{\text{ATE}}_{\cdot,j}^{\text{IPW}} - \text{ATE}_{\cdot,j}|$  and  $|\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}} - \text{ATE}_{\cdot,j}|$  across the  $Q$  simulation instances, respectively. We set  $r_p = r_\theta = \lfloor N^{1/5} \rfloor$  and vary  $N \in \{250, 500, 750, 1000, 1250, 1500\}$ . To make the scaling clear, we use least squares to produce the best  $N^{-\rho}$  fit to the maximum bias as  $N$  varies. We state the empirical decay rates in the legend, e.g., for DR, we report an empirical rate of  $N^{-0.62}$ . The DR estimator consistently outperforms the OI and IPW estimators.

## 7. Conclusion

This article introduces a new framework to estimate treatment effects in the presence unobserved confounding. We consider modern data-rich environments, where there are many units, and outcomes of interest per unit. We show it is possible to control for the confounding effects of a set of latent variables when this set is low-dimensional relative to the number of observed treatments and outcomes.

Our proposed estimator is doubly-robust, combining outcome imputation and inverse probability weighting with matrix completion. Analytical tractability of its distribution is gained through a novel cross-fitting procedure for matrix completion to estimate the treatment assignment probabilities and mean potential outcomes. We study the properties of the doubly-robust estimator, along with the outcome imputation and inverse probability weighting-based estimators under black-box matrix completion error rates. We show that the decay rate of the mean absolute error for the doubly-robust estimator dominates those of the outcome imputation and the inverse probability weighting estimators. Moreover, we establish a Gaussian approximation to the distribution of the doubly-robust estimator. Simulation results demonstrate the practical relevance of the formal properties of the doubly-robust estimator.

## Appendices

### A. Supporting Concentration and Convergence Results

This section presents known concentration bounds on subGaussian and subExponential random variables, along with the matrix Hoeffding bound and concludes with a basic result on convergence of random variables.

We use  $\text{subGaussian}(\sigma)$  to represent a subGaussian random variable, where  $\sigma$  is a bound on the subGaussian norm; and  $\text{subExponential}(\sigma)$  to represent a subExponential random variable, where  $\sigma$  is a bound on the subExponential norm. (Recall the definitions of the norms from Section 1.)

**Lemma A.1** (subGaussian concentration: Theorem 2.6.3 of Vershynin (2018)). *Let  $x \in \mathbb{R}^n$  be a random vector whose entries are independent, zero-mean, subGaussian( $\sigma$ ) random variables. Then, for any  $b \in \mathbb{R}^n$  and  $t \geq 0$ ,*

$$\mathbb{P}\left\{|b^\top x| \geq t\right\} \leq 2 \exp\left(\frac{-ct^2}{\sigma^2 \|b\|_2^2}\right).$$

The following corollary expresses the bound in Lemma A.1 in a convenient form.

**Corollary A.1** (subGaussian concentration). *Let  $x \in \mathbb{R}^n$  be a random vector whose entries are independent, zero-mean, subGaussian( $\sigma$ ) random variables. Then, for any  $b \in \mathbb{R}^n$  and any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,*

$$|b^\top x| \leq \sigma \sqrt{c\ell_\delta} \cdot \|b\|_2.$$

*Proof.* The proof follows from Lemma A.1 by choosing  $\delta \triangleq 2 \exp(-ct^2/\sigma^2 \|b\|_2^2)$ . □

**Lemma A.2** (subExponential concentration: Theorem 2.8.2 of Vershynin (2018)). *Let  $x \in \mathbb{R}^n$  be a random vector whose entries are independent, zero-mean, subExponential( $\sigma$ ) random variables. Then, for any  $b \in \mathbb{R}^n$  and  $t \geq 0$ ,*

$$\mathbb{P}\left\{|b^\top x| \geq t\right\} \leq 2 \exp\left(-c \min\left(\frac{t^2}{\sigma^2 \|b\|_2^2}, \frac{t}{\sigma \|b\|_\infty}\right)\right).$$

The following corollary expresses the bound in Lemma A.2 in a convenient form.

**Corollary A.2** (subExponential concentration). *Let  $x \in \mathbb{R}^n$  be a random vector whose entries are independent, zero-mean, subExponential( $\sigma$ ) random variables. Then, for any  $b \in \mathbb{R}^n$  and any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,*

$$|b^\top x| \leq \sigma m(c\ell_\delta) \cdot \|b\|_2,$$

where recall that  $m(c\ell_\delta) = \max(c\ell_\delta, \sqrt{c\ell_\delta})$ .

*Proof.* Choosing  $t = t_0 \sigma \|b\|_2$  in Lemma A.2, we have

$$\begin{aligned} \mathbb{P}\left\{|b^\top x| \geq t_0 \sigma \|b\|_2\right\} &\leq 2 \exp\left(-ct_0 \min\left(t_0, \frac{\|b\|_2}{\|b\|_\infty}\right)\right) \\ &\leq 2 \exp\left(-ct_0 \min(t_0, 1)\right), \end{aligned}$$

where the second inequality follows from  $\min\{t_0, c\} \geq \min\{t_0, 1\}$  for any  $c \geq 1$  and  $\|b\|_2 \geq \|b\|_\infty$ . Then, the proof follows by choosing  $\delta \triangleq 2 \exp(-ct_0 \min(t_0, 1))$  which fixes  $t_0 = \max\{\sqrt{c\ell_\delta}, c\ell_\delta\} = m(c\ell_\delta)$ .  $\square$

**Lemma A.3** (Product of subGaussians is subExponential: Lemma. 2.7.7 of Vershynin (2018)). *Let  $x_1$  and  $x_2$  be subGaussian( $\sigma_1$ ) and subGaussian( $\sigma_2$ ) random variables, respectively. Then,  $x_1 x_2$  is subExponential( $\sigma_1 \sigma_2$ ) random variable.*

**Lemma A.4** (Matrix Hoeffding bound: Theorem 1.3 of Tropp (2012)). *Let  $X_1, \dots, X_n$  be a sequence of independent, random, and symmetric matrices such that, for every  $i \in [N]$ ,  $X_i \in \mathbb{R}^{d \times d}$  and  $\mathbb{E}[X_i] = 0$ . Let  $A_1, \dots, A_n$  be a sequence of fixed symmetric matrices such that, for every  $i \in [N]$ ,  $A_i \in \mathbb{R}^{d \times d}$  and  $A_i^2 - X_i^2$  is positive semi-definite. Then,*

$$\mathbb{P}\left\{\lambda_{\max}\left(\sum_{i \in [n]} X_i\right) \geq t\right\} \leq d \exp\left(-\frac{t^2}{8\lambda_{\max}\left(\sum_{i \in [n]} A_i^2\right)}\right).$$

In the following corollary, we re-express the bound in Lemma A.4 in a convenient form.

**Corollary A.3** (Matrix Hoeffding bound). *Let  $X_1, \dots, X_n$  be a sequence of independent, random, and symmetric matrices such that, for every  $i \in [N]$ ,  $X_i \in \mathbb{R}^{d \times d}$  and  $\mathbb{E}[X_i] = 0$ . Let  $A_1, \dots, A_n$  be a sequence of fixed symmetric matrices such that, for every  $i \in [N]$ ,  $A_i \in \mathbb{R}^{d \times d}$  and  $A_i^2 - X_i^2$  is positive semi-definite. Then, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,*

$$\lambda_{\max}\left(\sum_{i \in [n]} X_i\right) \leq \sqrt{c\ell_{2\delta/d} \cdot \lambda_{\max}\left(\sum_{i \in [n]} A_i^2\right)}.$$

*Proof.* The proof follows from Lemma A.4 by choosing  $\delta \triangleq d \exp(-t^2/8 \cdot \lambda_{\max}(\sum_{i \in [n]} A_i^2))$ .  $\square$

Next lemma provides a useful intermediate result on convergence in probability.

**Lemma A.5.** *Let  $X_n$  and  $\overline{X}_n$  be sequences of random variables such that  $X_n = o_p(1)$ . Let  $\delta_n = o(1)$  be a deterministic sequence such that  $0 \leq \delta_n \leq 1$ . Suppose  $\mathbb{P}(|\overline{X}_n| \leq X_n) \geq 1 - \delta_n$ . Then,  $\overline{X}_n = o_p(1)$ .*

*Proof.* Consider any  $\epsilon > 0$ . Then, the event  $\{|\overline{X}_n| > \epsilon\}$  belongs to the union of  $\{|\overline{X}_n| > X_n\}$  and  $\{X_n > \epsilon\}$ . Using the union bound,

$$\mathbb{P}(|\overline{X}_n| > \epsilon) \leq \mathbb{P}(|\overline{X}_n| > X_n) + \mathbb{P}(X_n > \epsilon) \leq \delta_n + \mathbb{P}(X_n > \epsilon).$$

Then,  $\overline{X}_n = o_p(1)$  follows because  $X_n = o_p(1)$ .  $\square$



## B. Proof of Theorem 1: Finite Sample Guarantees for DR

Fix any  $j \in [M]$ . Recall the definitions Eqs. (5) and (11) of the parameters  $\text{ATE}_{\cdot,j}$  and corresponding doubly robust estimates  $\widehat{\text{ATE}}_{\cdot,j}^{\text{DR}}$ . The error  $\Delta \text{ATE}_{\cdot,j}^{\text{DR}} = \widehat{\text{ATE}}_{\cdot,j}^{\text{DR}} - \text{ATE}_{\cdot,j}$  can be re-expressed as

$$\begin{aligned} \Delta \text{ATE}_{\cdot,j}^{\text{DR}} &= \frac{1}{N} \sum_{i \in [N]} \left( \widehat{\theta}_{i,j}^{(1,\text{DR})} - \widehat{\theta}_{i,j}^{(0,\text{DR})} \right) - \frac{1}{N} \sum_{i \in [N]} \left( \theta_{i,j}^{(1)} - \theta_{i,j}^{(0)} \right) \\ &= \frac{1}{N} \sum_{i \in [N]} \left( \left( \widehat{\theta}_{i,j}^{(1,\text{DR})} - \theta_{i,j}^{(1)} \right) - \left( \widehat{\theta}_{i,j}^{(0,\text{DR})} - \theta_{i,j}^{(0)} \right) \right) \\ &\stackrel{(a)}{=} \frac{1}{N} \sum_{i \in [N]} \left( \mathbb{T}_{i,j}^{(1,\text{DR})} + \mathbb{T}_{i,j}^{(0,\text{DR})} \right), \end{aligned} \quad (\text{A.1})$$

where (a) follows after defining  $\mathbb{T}_{i,j}^{(1,\text{DR})} \triangleq \left( \widehat{\theta}_{i,j}^{(1,\text{DR})} - \theta_{i,j}^{(1)} \right)$  and  $\mathbb{T}_{i,j}^{(0,\text{DR})} \triangleq -\left( \widehat{\theta}_{i,j}^{(0,\text{DR})} - \theta_{i,j}^{(0)} \right)$  for every  $(i, j) \in [N] \times [M]$ . Then, we have

$$\begin{aligned} \mathbb{T}_{i,j}^{(1,\text{DR})} &= \widehat{\theta}_{i,j}^{(1,\text{DR})} - \theta_{i,j}^{(1)} \\ &\stackrel{(a)}{=} \widehat{\theta}_{i,j}^{(1)} + (y_{i,j} - \widehat{\theta}_{i,j}^{(1)}) \frac{a_{i,j}}{\widehat{p}_{i,j}} - \theta_{i,j}^{(1)} \\ &\stackrel{(b)}{=} \widehat{\theta}_{i,j}^{(1)} + (\theta_{i,j}^{(1)} + \varepsilon_{i,j}^{(1)} - \widehat{\theta}_{i,j}^{(1)}) \frac{p_{i,j} + \eta_{i,j}}{\widehat{p}_{i,j}} - \theta_{i,j}^{(1)} \\ &= (\widehat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) \left( 1 - \frac{p_{i,j} + \eta_{i,j}}{\widehat{p}_{i,j}} \right) + \varepsilon_{i,j}^{(1)} \left( \frac{p_{i,j} + \eta_{i,j}}{\widehat{p}_{i,j}} \right) \\ &= \frac{(\widehat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) (\widehat{p}_{i,j} - p_{i,j})}{\widehat{p}_{i,j}} - \frac{(\widehat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) \eta_{i,j}}{\widehat{p}_{i,j}} + \frac{\varepsilon_{i,j}^{(1)} p_{i,j}}{\widehat{p}_{i,j}} + \frac{\varepsilon_{i,j}^{(1)} \eta_{i,j}}{\widehat{p}_{i,j}}, \end{aligned} \quad (\text{A.2})$$

where (a) follows from Eq. (12), and (b) follows from Eqs. (1) to (3). A similar derivation for  $a = 0$  implies that

$$\begin{aligned} \mathbb{T}_{i,j}^{(0,\text{DR})} &= -\widehat{\theta}_{i,j}^{(0,\text{DR})} + \theta_{i,j}^{(0)} \\ &= -\frac{(\widehat{\theta}_{i,j}^{(0)} - \theta_{i,j}^{(0)}) (1 - \widehat{p}_{i,j} - (1 - p_{i,j}))}{1 - \widehat{p}_{i,j}} + \frac{(\widehat{\theta}_{i,j}^{(0)} - \theta_{i,j}^{(0)}) (-\eta_{i,j})}{1 - \widehat{p}_{i,j}} - \frac{\varepsilon_{i,j}^{(0)} (1 - p_{i,j})}{1 - \widehat{p}_{i,j}} - \frac{\varepsilon_{i,j}^{(0)} (-\eta_{i,j})}{1 - \widehat{p}_{i,j}} \\ &= \frac{(\widehat{\theta}_{i,j}^{(0)} - \theta_{i,j}^{(0)}) (\widehat{p}_{i,j} - p_{i,j})}{1 - \widehat{p}_{i,j}} - \frac{(\widehat{\theta}_{i,j}^{(0)} - \theta_{i,j}^{(0)}) \eta_{i,j}}{1 - \widehat{p}_{i,j}} - \frac{\varepsilon_{i,j}^{(0)} (1 - p_{i,j})}{1 - \widehat{p}_{i,j}} + \frac{\varepsilon_{i,j}^{(0)} \eta_{i,j}}{1 - \widehat{p}_{i,j}}. \end{aligned} \quad (\text{A.3})$$

Consider any  $a \in \{0, 1\}$  and any  $\delta \in (0, 1)$ . We claim that, with probability at least  $1 - 6\delta$ ,

$$\frac{1}{N} \left| \sum_{i \in [N]} \mathbb{T}_{i,j}^{(a,\text{DR})} \right| \leq \frac{2}{\lambda} \mathcal{E}(\widehat{\Theta}^{(a)}) \cdot \mathcal{E}(\widehat{P}) + \frac{2\sqrt{c\ell_\delta}}{\lambda\sqrt{\ell_1 N}} \mathcal{E}(\widehat{\Theta}^{(a)}) + \frac{2\bar{\sigma}\sqrt{c\ell_\delta}}{\lambda\sqrt{N}} + \frac{2\bar{\sigma}m(c\ell_\delta)}{\lambda\sqrt{\ell_1 N}}, \quad (\text{A.4})$$

where recall that  $m(cl_\delta) = \max(cl_\delta, \sqrt{cl_\delta})$ . We provide a proof of this claim at the end of this section. Applying triangle inequality in Eq. (A.1) and using Eq. (A.4) with a union bound, we obtain that

$$|\Delta \text{ATE}_{\cdot,j}^{\text{DR}}| \leq \frac{2}{\lambda} \cdot \mathcal{E}(\hat{\Theta}) \mathcal{E}(\hat{P}) + \frac{2\sqrt{cl_\delta}}{\lambda\sqrt{\ell_1 N}} \mathcal{E}(\hat{\Theta}) + \frac{4\bar{\sigma}\sqrt{cl_\delta}}{\lambda\sqrt{N}} + \frac{4\bar{\sigma}m(cl_\delta)}{\lambda\sqrt{\ell_1 N}}, \quad (\text{A.5})$$

with probability at least  $1 - 12\delta$ . The claim in Eq. (18) follows by re-parameterizing  $\delta$ .

Next, to establish the claim in Eq. (19), choose  $\delta = 1/N$  and note that every term in the right hand side of Eq. (A.5) is  $o_p(1)$  under the conditions on  $\mathcal{E}(\hat{\Theta})$  and  $\mathcal{E}(\hat{P})$ . Then, Eq. (19) follows from Lemma A.5.

**Proof of bound (A.4).** Recall the partitioning of the units  $[N]$  into  $\mathcal{R}_0$  and  $\mathcal{R}_1$  from Assumption 4. Condition on this partition. Now, to enable the application of concentration bounds, we split the summation over  $i \in [N]$  in the left hand side of Eq. (A.4) into two parts—one over  $i \in \mathcal{R}_0$  and the other over  $i \in \mathcal{R}_1$ —such that the noise terms are independent of the estimates of  $\Theta^{(0)}, \Theta^{(1)}, P$  in each of these parts as in Eqs. (14) and (15).

Note that  $|\sum_{i \in [N]} \mathbb{T}_{i,j}^{(1,\text{DR})}| \leq |\sum_{i \in \mathcal{R}_0} \mathbb{T}_{i,j}^{(1,\text{DR})}| + |\sum_{i \in \mathcal{R}_1} \mathbb{T}_{i,j}^{(1,\text{DR})}|$ . Let  $s \in \{0, 1\}$ . Eq. (A.2) and triangle inequality imply

$$\begin{aligned} \left| \sum_{i \in \mathcal{R}_s} \mathbb{T}_{i,j}^{(1,\text{DR})} \right| &\leq \left| \sum_{i \in \mathcal{R}_s} \frac{(\hat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) (\hat{p}_{i,j} - p_{i,j})}{\hat{p}_{i,j}} \right| + \left| \sum_{i \in \mathcal{R}_s} \frac{(\hat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) \eta_{i,j}}{\hat{p}_{i,j}} \right| \\ &\quad + \left| \sum_{i \in \mathcal{R}_s} \frac{\varepsilon_{i,j}^{(1)} p_{i,j}}{\hat{p}_{i,j}} \right| + \left| \sum_{i \in \mathcal{R}_s} \frac{\varepsilon_{i,j}^{(1)} \eta_{i,j}}{\hat{p}_{i,j}} \right|. \end{aligned} \quad (\text{A.6})$$

Applying the Cauchy-Schwarz inequality to bound the first term yields that

$$\begin{aligned} \left| \sum_{i \in \mathcal{R}_s} \frac{(\hat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) (\hat{p}_{i,j} - p_{i,j})}{\hat{p}_{i,j}} \right| &\leq \sqrt{\sum_{i \in \mathcal{R}_s} \left( \frac{\hat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}}{\hat{p}_{i,j}} \right)^2 \cdot \sum_{i \in \mathcal{R}_s} (\hat{p}_{i,j} - p_{i,j})^2} \\ &\leq \|(\hat{\Theta}_{\cdot,j}^{(1)} - \Theta_{\cdot,j}^{(1)}) \odot \hat{P}_{\cdot,j}\|_2 \| \hat{P}_{\cdot,j} - P_{\cdot,j} \|_2. \end{aligned} \quad (\text{A.7})$$

To bound the second term in Eq. (A.6), note that  $\eta_{i,j}$  is  $\text{subGaussian}(1/\sqrt{\ell_1})$  (see Example 2.5.8 in Vershynin (2018)), zero-mean due to Assumption 2(a), and independent across all  $i \in [N]$  due to Assumption 2(c). Moreover, Assumption 4 (i.e., Eq. (14)) provides that  $(\hat{\theta}_{i,j}^{(1)}, \hat{\theta}_{i,j}^{(0)}, \hat{p}_{i,j})_{i \in \mathcal{R}_s} \perp\!\!\!\perp (\eta_{i,j})_{i \in \mathcal{R}_s}$ . Hence, applying the subGaussian concentration (Corollary A.1) for  $(\eta_{i,j})_{i \in \mathcal{R}_s}$  yields that

$$\left| \sum_{i \in \mathcal{R}_s} \frac{(\hat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) \eta_{i,j}}{\hat{p}_{i,j}} \right| \leq \frac{\sqrt{cl_\delta}}{\sqrt{\ell_1}} \cdot \sqrt{\sum_{i \in \mathcal{R}_s} \left( \frac{\hat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}}{\hat{p}_{i,j}} \right)^2} \leq \frac{\sqrt{cl_\delta}}{\sqrt{\ell_1}} \|(\hat{\Theta}_{\cdot,j}^{(1)} - \Theta_{\cdot,j}^{(1)}) \odot \hat{P}_{\cdot,j}\|_2, \quad (\text{A.8})$$

with probability at least  $1 - \delta$ .

To bound the third term in Eq. (A.6), note that  $\varepsilon_{i,j}^{(1)}$  is subGaussian( $\bar{\sigma}$ ) due to Assumption 2(e), zero-mean due to Assumption 2(a), and independent across all  $i \in [N]$  due to Assumption 2(d). Moreover, Assumption 4 provides (i.e., Eq. (15)) that  $(\hat{p}_{i,j})_{i \in \mathcal{R}_s} \perp\!\!\!\perp (\varepsilon_{i,j}^{(1)})_{i \in \mathcal{R}_s}$ . Hence, applying the subGaussian concentration (Corollary A.1) for  $(\varepsilon_{i,j}^{(1)})_{i \in \mathcal{R}_s}$  yields that

$$\left| \sum_{i \in \mathcal{R}_s} \frac{\varepsilon_{i,j}^{(1)} p_{i,j}}{\hat{p}_{i,j}} \right| \leq \bar{\sigma} \sqrt{c\ell_\delta} \sqrt{\sum_{i \in \mathcal{R}_s} \left( \frac{p_{i,j}}{\hat{p}_{i,j}} \right)^2} \leq \bar{\sigma} \sqrt{c\ell_\delta} \|P_{\cdot,j} \odot \hat{P}_{\cdot,j}\|_2, \quad (\text{A.9})$$

with probability at least  $1 - \delta$ .

Finally, to bound the fourth term in Eq. (A.6), note that  $\varepsilon_{i,j}^{(1)} \eta_{i,j}$  is subExponential( $\bar{\sigma}/\sqrt{\ell_1}$ ) due to Lemma A.3. Further,  $\varepsilon_{i,j}^{(1)} \eta_{i,j}$  is zero-mean due to Assumption 2(a) and independent across all  $i \in [N]$  due to Assumption 2(b) to (d). Moreover, Assumption 4 (i.e., Eqs. (14) and (15)) imply that  $(\hat{p}_{i,j})_{i \in \mathcal{R}_s} \perp\!\!\!\perp (\eta_{i,j}, \varepsilon_{i,j}^{(1)})_{i \in \mathcal{R}_s}$ . Hence, applying the subExponential concentration (Corollary A.2) for  $(\eta_{i,j} \varepsilon_{i,j}^{(1)})_{i \in \mathcal{R}_s}$  yields that

$$\left| \sum_{i \in \mathcal{R}_s} \frac{\varepsilon_{i,j}^{(1)} \eta_{i,j}}{\hat{p}_{i,j}} \right| \leq \frac{\bar{\sigma} m(c\ell_\delta)}{\sqrt{\ell_1}} \|\mathbf{1}_N \odot \hat{P}_{\cdot,j}\|_2, \quad (\text{A.10})$$

with probability at least  $1 - \delta$ . Putting together Eqs. (A.6) to (A.10), we conclude that, with probability at least  $1 - 3\delta$ ,

$$\begin{aligned} \frac{1}{N} \left| \sum_{i \in \mathcal{R}_s} \mathbb{T}_{i,j}^{(1, \text{DR})} \right| &\leq \frac{1}{N} \|(\hat{\Theta}_{\cdot,j}^{(1)} - \Theta_{\cdot,j}^{(1)}) \odot \hat{P}_{\cdot,j}\|_2 \|\hat{P}_{\cdot,j} - P_{\cdot,j}\|_2 + \frac{\sqrt{c\ell_\delta}}{\sqrt{\ell_1} N} \|(\hat{\Theta}_{\cdot,j}^{(1)} - \Theta_{\cdot,j}^{(1)}) \odot \hat{P}_{\cdot,j}\|_2 \\ &\quad + \frac{\bar{\sigma} \sqrt{c\ell_\delta}}{N} \|P_{\cdot,j} \odot \hat{P}_{\cdot,j}\|_2 + \frac{\bar{\sigma} m(c\ell_\delta)}{\sqrt{\ell_1} N} \|\mathbf{1}_N \odot \hat{P}_{\cdot,j}\|_2. \end{aligned} \quad (\text{A.11})$$

Then, noting that  $1/\hat{p}_{i,j} \leq 1/\bar{\lambda}$  for every  $i \in [N]$  and  $j \in [M]$  from Assumption 3, and consequently that  $\|B_{\cdot,j} \odot \hat{P}_{\cdot,j}\|_2 \leq \|B\|_{1,2}/\bar{\lambda}$  for any matrix  $B$  and every  $j \in [M]$ , we obtain the following bound, with probability at least  $1 - 3\delta$ ,

$$\begin{aligned} \frac{1}{N} \left| \sum_{i \in \mathcal{R}_s} \mathbb{T}_{i,j}^{(1, \text{DR})} \right| &\leq \frac{1}{\bar{\lambda} N} \|\hat{\Theta}^{(1)} - \Theta^{(1)}\|_{1,2} \|\hat{P} - P\|_{1,2} + \frac{\sqrt{c\ell_\delta}}{\bar{\lambda} \sqrt{\ell_1} N} \|\hat{\Theta}^{(1)} - \Theta^{(1)}\|_{1,2} \\ &\quad + \frac{\bar{\sigma} \sqrt{c\ell_\delta}}{\bar{\lambda} N} \|P\|_{1,2} + \frac{\bar{\sigma} m(c\ell_\delta)}{\bar{\lambda} \sqrt{\ell_1} N} \|\mathbf{1}\|_{1,2} \end{aligned} \quad (\text{A.12})$$

$$\stackrel{(a)}{\leq} \frac{1}{\bar{\lambda}} \mathcal{E}(\hat{\Theta}^{(1)}) \cdot \mathcal{E}(\hat{P}) + \frac{\sqrt{c\ell_\delta}}{\bar{\lambda} \sqrt{\ell_1} N} \mathcal{E}(\hat{\Theta}^{(1)}) + \frac{\bar{\sigma} \sqrt{c\ell_\delta}}{\bar{\lambda} \sqrt{N}} + \frac{\bar{\sigma} m(c\ell_\delta)}{\bar{\lambda} \sqrt{\ell_1} N}, \quad (\text{A.13})$$

where (a) follows from Eq. (16) and because  $\|P\|_{1,2} \leq \sqrt{N}$  and  $\|\mathbf{1}\|_{1,2} = \sqrt{N}$ . Then, the claim in Eq. (A.4) follows for  $a = 1$  by using Eq. (A.13) and applying a union bound over  $s \in \{0, 1\}$ . The proof of Eq. (A.4) for  $a = 0$  follows similarly.

## C. Proof of Theorem 2: Asymptotic Normality for DR

For every  $(i, j) \in [N] \times [M]$ , recall the definitions of  $\mathbb{T}_{i,j}^{(1, \text{DR})}$  and  $\mathbb{T}_{i,j}^{(0, \text{DR})}$  from Eq. (A.2) and Eq. (A.3), respectively. Then, define

$$\begin{aligned}\mathbb{X}_{i,j}^{(1, \text{DR})} &\triangleq \mathbb{T}_{i,j}^{(1, \text{DR})} - \varepsilon_{i,j}^{(1)} - \frac{\varepsilon_{i,j}^{(1)} \eta_{i,j}}{p_{i,j}} \\ \mathbb{X}_{i,j}^{(0, \text{DR})} &\triangleq \mathbb{T}_{i,j}^{(0, \text{DR})} + \varepsilon_{i,j}^{(0)} - \frac{\varepsilon_{i,j}^{(0)} \eta_{i,j}}{1 - p_{i,j}},\end{aligned}\tag{A.14}$$

and

$$\mathbb{Z}_{i,j}^{\text{DR}} \triangleq \varepsilon_{i,j}^{(1)} + \frac{\varepsilon_{i,j}^{(1)} \eta_{i,j}}{p_{i,j}} - \varepsilon_{i,j}^{(0)} + \frac{\varepsilon_{i,j}^{(0)} \eta_{i,j}}{1 - p_{i,j}}.\tag{A.15}$$

Fix any  $j \in [M]$ . Then, the simplification of  $\Delta \text{ATE}_{\cdot,j}^{\text{DR}}$  in Eq. (A.1) can be re-expressed as

$$\Delta \text{ATE}_{\cdot,j}^{\text{DR}} = \frac{1}{N} \sum_{i \in [N]} \left( \mathbb{X}_{i,j}^{(1, \text{DR})} + \mathbb{X}_{i,j}^{(0, \text{DR})} + \mathbb{Z}_{i,j}^{\text{DR}} \right)\tag{A.16}$$

We prove in Appendices C.1 and C.2 the following convergence results for the above terms.

**Lemma C.1** (Convergence of  $\mathbb{X}_j^{\text{DR}}$ ). *Suppose Assumptions 1 to 4 and conditions (C1) to (C3) in Theorem 2 hold. For any fixed  $j \in [M]$ ,*

$$\frac{1}{\bar{\sigma}_j \sqrt{N}} \sum_{i \in [N]} \left( \mathbb{X}_{i,j}^{(1, \text{DR})} + \mathbb{X}_{i,j}^{(0, \text{DR})} \right) = o_p(1).$$

**Lemma C.2** (Convergence of  $\mathbb{Z}_j^{\text{DR}}$ ). *Suppose Assumptions 1 and 2 hold and condition (C3) in Theorem 2 hold. For any fixed*

$$\frac{1}{\bar{\sigma}_j \sqrt{N}} \sum_{i \in [N]} \mathbb{Z}_{i,j}^{\text{DR}} \xrightarrow{d} \mathcal{N}(0, 1).$$

Now, Theorem 2 follows by applying Slutsky's theorem to put together Lemmas C.1 and C.2 with Eq. (A.16).

### C.1. Proof of Lemma C.1

Fix any  $j \in [M]$ . Consider any  $a \in \{0, 1\}$  and any  $\delta \in (0, 1)$ . We claim that, with probability at least  $1 - \delta/2$ ,

$$\frac{1}{N} \sum_{i \in [N]} \mathbb{X}_{i,j}^{(a, \text{DR})} \leq \frac{2}{\bar{\lambda}} \cdot \mathcal{E}(\hat{\Theta}^{(a)}) \mathcal{E}(\hat{P}) + \frac{2\sqrt{c\ell_{\delta/12}}}{\bar{\lambda}\sqrt{\ell_1}} \cdot \frac{\mathcal{E}(\hat{\Theta}^{(a)})}{\sqrt{N}} + \frac{4\bar{\sigma}m(c\ell_{\delta/12})}{\lambda\bar{\lambda}\sqrt{\ell_1}} \cdot \frac{\mathcal{E}(\hat{P})}{\sqrt{N}},\tag{A.17}$$

where recall that  $m(c\ell_{\delta/12}) = \max(c\ell_{\delta/12}, \sqrt{c\ell_{\delta/12}})$ . We provide a proof of this claim at the end of this section. Then, using Eq. (A.17) with a union bound, and the fact that  $\bar{\sigma}_j \geq c > 0$  as per condition (C3), we obtain the following with probability at least  $1 - \delta$ ,

$$\frac{1}{\bar{\sigma}_j \sqrt{N}} \sum_{\substack{i \in [N], \\ a \in \{0,1\}}} \mathbb{X}_{i,j}^{(a,\text{DR})} \leq \frac{1}{c} \left( \frac{2}{\lambda} \cdot \sqrt{N} \mathcal{E}(\hat{\Theta}) \mathcal{E}(\hat{P}) + \frac{2\sqrt{c\ell_{\delta/12}}}{\lambda\sqrt{\ell_1}} \cdot \mathcal{E}(\hat{\Theta}) + \frac{8\bar{\sigma}m(c\ell_{\delta/12})}{\lambda\lambda\sqrt{\ell_1}} \cdot \mathcal{E}(\hat{P}) \right). \quad (\text{A.18})$$

We emphasize that Eq. (A.18) holds for any  $\delta \in (0, 1)$ . Next, we choose a particular  $\delta$  that is  $o(1)$  and, under conditions (C1) and (C2), show that each of the three terms in the right hand side of Eq. (A.18) are  $o_p(1)$ . In particular, we choose

$$\delta = \exp \left( -1 / \max \{t_N, \sqrt{s_N}\} \right).$$

We note that this choice of  $\delta$  suffices. First,  $\delta = o(1)$  follows by using condition (C1), the continuous mapping theorem, and the convergence in probability of the maximum of two sequences of variables. Second,  $\sqrt{N} \mathcal{E}(\hat{\Theta}) \mathcal{E}(\hat{P}) = o_p(1)$  from condition (C2). Third,  $\sqrt{\ell_{\delta/12}} \mathcal{E}(\hat{\Theta}) = o_p(1)$  follows by using condition (C1) and the continuous mapping theorem after noting that  $\sqrt{\ell_{\delta/12}} \mathcal{E}(\hat{\Theta}) \leq O_p(t_N^{1/2})$ . Fourth,  $m(\ell_{\delta/12}) \mathcal{E}(\hat{P}) = o_p(1)$  follows by using condition (C1) and the continuous mapping theorem after noting that  $m(\ell_{\delta/12}) \mathcal{E}(\hat{P}) \leq O_p(\max\{s_N^{1/2}, s_N^{3/4}\})$ . Finally, Lemma C.1 follows from Lemma A.5.

**Proof of Eq. (A.17)** This proof follows a very similar road map to that used for establishing the inequality in display (A.4). Recall the partitioning of the units  $[N]$  into  $\mathcal{R}_0$  and  $\mathcal{R}_1$  from Assumption 4. Condition on this partition. Now, to enable the application of concentration bounds, we split the summation over  $i \in [N]$  in the left hand side of Eq. (A.17) into two parts—one over  $i \in \mathcal{R}_0$  and the other over  $i \in \mathcal{R}_1$ —such that the noise terms are independent of the estimates of  $\Theta^{(0)}, \Theta^{(1)}, P$  in each of these parts as in Eqs. (14) and (15).

Fix  $a = 1$ . Then, Eqs. (A.2) and (A.14) imply that

$$\begin{aligned} \mathbb{X}_{i,j}^{(1,\text{DR})} &= \frac{(\hat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) (\hat{p}_{i,j} - p_{i,j})}{\hat{p}_{i,j}} - \frac{(\hat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) \eta_{i,j}}{\hat{p}_{i,j}} + \frac{\varepsilon_{i,j}^{(1)} p_{i,j}}{\hat{p}_{i,j}} + \frac{\varepsilon_{i,j}^{(1)} \eta_{i,j}}{\hat{p}_{i,j}} - \varepsilon_{i,j}^{(1)} - \frac{\varepsilon_{i,j}^{(1)} \eta_{i,j}}{p_{i,j}} \\ &= \frac{(\hat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) (\hat{p}_{i,j} - p_{i,j})}{\hat{p}_{i,j}} - \frac{(\hat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) \eta_{i,j}}{\hat{p}_{i,j}} - \frac{\varepsilon_{i,j}^{(1)} (\hat{p}_{i,j} - p_{i,j})}{\hat{p}_{i,j}} - \frac{\varepsilon_{i,j}^{(1)} \eta_{i,j} (\hat{p}_{i,j} - p_{i,j})}{\hat{p}_{i,j} p_{i,j}}. \end{aligned}$$

Now, note that  $|\sum_{i \in [N]} \mathbb{X}_{i,j}^{(1,\text{DR})}| \leq |\sum_{i \in \mathcal{R}_0} \mathbb{X}_{i,j}^{(1,\text{DR})}| + |\sum_{i \in \mathcal{R}_1} \mathbb{X}_{i,j}^{(1,\text{DR})}|$ . Fix any  $s \in \{0, 1\}$ . Then, triangle inequality implies that

$$\begin{aligned} \left| \sum_{i \in \mathcal{R}_s} \mathbb{X}_{i,j}^{(1,\text{DR})} \right| &\leq \left| \sum_{i \in \mathcal{R}_s} \frac{(\hat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) (\hat{p}_{i,j} - p_{i,j})}{\hat{p}_{i,j}} \right| + \left| \sum_{i \in \mathcal{R}_s} \frac{(\hat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) \eta_{i,j}}{\hat{p}_{i,j}} \right| \\ &\quad + \left| \sum_{i \in \mathcal{R}_s} \frac{\varepsilon_{i,j}^{(1)} (\hat{p}_{i,j} - p_{i,j})}{\hat{p}_{i,j}} \right| + \left| \sum_{i \in \mathcal{R}_s} \frac{\varepsilon_{i,j}^{(1)} \eta_{i,j} (\hat{p}_{i,j} - p_{i,j})}{\hat{p}_{i,j} p_{i,j}} \right|. \quad (\text{A.19}) \end{aligned}$$

Next, note that the decomposition in Eq. (A.19) is identical to the one in Eq. (A.6), except for the fact when compared to Eq. (A.6), the last two terms in Eq. (A.19) have an additional factor of  $(\hat{p}_{i,j} - p_{i,j})/p_{i,j}$ . As a result, mimicking steps used to derive Eq. (A.11), we can obtain the following bound, with probability at least  $1 - 3\delta$ ,

$$\begin{aligned} \frac{1}{N} \left| \sum_{i \in \mathcal{R}_s} \mathbb{X}_{i,j}^{(1, \text{DR})} \right| &\leq \frac{1}{N} \|(\hat{\Theta}_{\cdot,j}^{(1)} - \Theta_{\cdot,j}^{(1)}) \odot \hat{P}_{\cdot,j}\|_2 \|\hat{P}_{\cdot,j} - P_{\cdot,j}\|_2 + \frac{\sqrt{cl_\delta}}{\sqrt{\ell_1}N} \|(\hat{\Theta}_{\cdot,j}^{(1)} - \Theta_{\cdot,j}^{(1)}) \odot \hat{P}_{\cdot,j}\|_2 \\ &\quad + \frac{\bar{\sigma}\sqrt{cl_\delta}}{N} \|(\hat{P}_{\cdot,j} - P_{\cdot,j}) \odot \hat{P}_{\cdot,j}\|_2 + \frac{\bar{\sigma}m(cl_\delta)}{\sqrt{\ell_1}N} \|(\hat{P}_{\cdot,j} - P_{\cdot,j}) \odot (\hat{P}_{\cdot,j} \odot P_{\cdot,j})\|_2. \end{aligned}$$

Then, noting that  $1/p_{i,j} \leq 1/\lambda$  and  $1/\hat{p}_{i,j} \leq 1/\bar{\lambda}$  for all  $i \in [N]$  and  $j \in [M]$  from Assumptions 1 and 3, and consequently that  $\|B_{\cdot,j} \odot \hat{P}_{\cdot,j}\|_2 \leq \|B\|_{1,2}/\bar{\lambda}$  and  $\|B_{\cdot,j} \odot P_{\cdot,j}\|_2 \leq \|B\|_{1,2}/\lambda$  for any matrix  $B$  and every  $j \in [M]$ , we obtain the following bound, with probability at least  $1 - 3\delta$ ,

$$\begin{aligned} \frac{1}{N} \left| \sum_{i \in \mathcal{R}_s} \mathbb{X}_{i,j}^{(1, \text{DR})} \right| &\leq \frac{1}{\bar{\lambda}N} \|\hat{\Theta}^{(1)} - \Theta^{(1)}\|_{1,2} \|\hat{P} - P\|_{1,2} + \frac{\sqrt{cl_\delta}}{\bar{\lambda}\sqrt{\ell_1}N} \|\hat{\Theta}^{(1)} - \Theta^{(1)}\|_{1,2} \\ &\quad + \frac{\bar{\sigma}\sqrt{cl_\delta}}{\bar{\lambda}N} \|\hat{P} - P\|_{1,2} + \frac{\bar{\sigma}m(cl_\delta)}{\lambda\bar{\lambda}\sqrt{\ell_1}N} \|\hat{P} - P\|_{1,2} \\ &\stackrel{(a)}{\leq} \frac{1}{\bar{\lambda}N} \|\hat{\Theta}^{(1)} - \Theta^{(1)}\|_{1,2} \|\hat{P} - P\|_{1,2} + \frac{\sqrt{cl_\delta}}{\bar{\lambda}\sqrt{\ell_1}N} \|\hat{\Theta}^{(1)} - \Theta^{(1)}\|_{1,2} \\ &\quad + \frac{2\bar{\sigma}m(cl_\delta)}{\lambda\bar{\lambda}\sqrt{\ell_1}N} \|\hat{P} - P\|_{1,2} \\ &\stackrel{(b)}{\leq} \frac{1}{\bar{\lambda}} \mathcal{E}(\hat{\Theta}^{(1)}) \cdot \mathcal{E}(\hat{P}) + \frac{\sqrt{cl_\delta}}{\bar{\lambda}\sqrt{\ell_1}N} \mathcal{E}(\hat{\Theta}^{(1)}) + \frac{2\bar{\sigma}m(cl_\delta)}{\lambda\bar{\lambda}\sqrt{\ell_1}N} \mathcal{E}(\hat{P}), \end{aligned} \tag{A.20}$$

where (a) follows because  $\lambda \leq 1/2 < 1/\sqrt{\ell_1}$  from Assumption 1 and  $\sqrt{cl_\delta} \leq m(cl_\delta)$ , and (b) follows from Eq. (16). Then, the claim in Eq. (A.17) follows for  $a = 1$  by applying a union bound over  $s \in \{0, 1\}$  using Eq. (A.20), and re-parameterizing  $\delta$ . The proof of Eq. (A.4) for  $a = 0$  follows similarly.

## C.2. Proof of Lemma C.2

To prove this result, we invoke Lyapunov central limit theorem (CLT).

**Lemma C.3** (Lyapunov CLT, see Theorem 27.3 of Billingsley (2017)). *Consider a sequence  $x_1, x_2, \dots$  of independent, mean-zero, and finite variance random variables. If Lyapunov's condition is satisfied, i.e., there exists  $\omega > 0$  such that*

$$\frac{\sum_{i=1}^N \mathbb{E}[|x_i|^{2+\omega}]}{(\sum_{i=1}^N \mathbb{E}[x_i^2])^{\frac{2+\omega}{2}}} \rightarrow 0, \tag{A.21}$$

as  $N \rightarrow \infty$ , then

$$\frac{\sum_{i=1}^N x_i}{(\sum_{i=1}^N \mathbb{E}[x_i^2])^{\frac{1}{2}}} \xrightarrow{d} \mathcal{N}(0, 1),$$

as  $N \rightarrow \infty$ .

Fix any  $j \in [M]$ . We apply Lyapunov CLT in Lemma C.3 on the sequence  $\mathbb{Z}_{1,j}^{\text{DR}}, \mathbb{Z}_{2,j}^{\text{DR}}, \dots$  where  $\mathbb{Z}_{i,j}^{\text{DR}}$  is as defined in Eq. (A.15). Note that Assumption 2(a) and (b) imply  $\mathbb{E}[\mathbb{Z}_{i,j}^{\text{DR}}] = 0$  for all  $i \in [N]$ , and Assumption 2(b) to (d) imply that  $\mathbb{Z}_{i,j}^{\text{DR}} \perp \mathbb{Z}_{i',j}^{\text{DR}}$  for all  $i \neq i' \in [N]$ . First, we show in Appendix C.2.1 that

$$\mathbb{V}\text{ar}(\mathbb{Z}_{i,j}^{\text{DR}}) = \frac{(\sigma_{i,j}^{(1)})^2}{p_{i,j}} + \frac{(\sigma_{i,j}^{(0)})^2}{1 - p_{i,j}}, \quad (\text{A.22})$$

for each  $i \in [N]$ . Next, we show in Appendix C.2.2 that Lyapunov's condition (A.21) holds for the sequence  $\mathbb{Z}_{1,j}^{\text{DR}}, \mathbb{Z}_{2,j}^{\text{DR}}, \dots$  with  $\omega = 1$ . Finally, applying Lemma C.3 and using the definition of  $\bar{\sigma}_j$  from Eq. (22) yields Lemma C.2.

#### C.2.1. Proof of Eq. (A.22)

Fix any  $i \in [N]$  and consider  $\mathbb{V}\text{ar}(\mathbb{Z}_{i,j}^{\text{DR}})$ . We have

$$\mathbb{V}\text{ar}(\mathbb{Z}_{i,j}^{\text{DR}}) = \mathbb{V}\text{ar}\left(\varepsilon_{i,j}^{(1)}\left(1 + \frac{\eta_{i,j}}{p_{i,j}}\right) - \varepsilon_{i,j}^{(0)}\left(1 - \frac{\eta_{i,j}}{1 - p_{i,j}}\right)\right). \quad (\text{A.23})$$

We claim the following:

$$\mathbb{V}\text{ar}\left(\varepsilon_{i,j}^{(1)}\left(1 + \frac{\eta_{i,j}}{p_{i,j}}\right)\right) = \frac{(\sigma_{i,j}^{(1)})^2}{p_{i,j}}, \quad (\text{A.24})$$

$$\mathbb{V}\text{ar}\left(\varepsilon_{i,j}^{(0)}\left(1 - \frac{\eta_{i,j}}{1 - p_{i,j}}\right)\right) = \frac{(\sigma_{i,j}^{(0)})^2}{1 - p_{i,j}}, \quad \text{and} \quad (\text{A.25})$$

$$\mathbb{C}\text{ov}\left(\varepsilon_{i,j}^{(1)}\left(1 + \frac{\eta_{i,j}}{p_{i,j}}\right), \varepsilon_{i,j}^{(0)}\left(1 - \frac{\eta_{i,j}}{1 - p_{i,j}}\right)\right) = 0. \quad (\text{A.26})$$

Then, Eq. (A.22) follows by putting together Eqs. (A.23) to (A.26) by using  $\mathbb{V}\text{ar}(x_1 - x_2) = \mathbb{V}\text{ar}(x_1) + \mathbb{V}\text{ar}(x_2) - 2\mathbb{C}\text{ov}(x_1, x_2)$  for any random variables  $x_1$  and  $x_2$ . It remains to establish the claims in Eqs. (A.24) to (A.26).

Assumption 2 immediately implies that  $\varepsilon_{i,j}^{(1)} \perp \eta_{i,j}$  and  $\mathbb{E}[\varepsilon_{i,j}^{(1)}] = \mathbb{E}[\eta_{i,j}] = 0$  (so that  $\varepsilon_{i,j}^{(1)}(1 + \frac{\eta_{i,j}}{p_{i,j}})$  is mean zero). Applying these observations, we obtain Eq. (A.24) as follows,

$$\begin{aligned} \mathbb{V}\text{ar}\left(\varepsilon_{i,j}^{(1)}\left(1 + \frac{\eta_{i,j}}{p_{i,j}}\right)\right) &= \mathbb{E}\left[\left(\varepsilon_{i,j}^{(1)}\left(1 + \frac{\eta_{i,j}}{p_{i,j}}\right)\right)^2\right] = \mathbb{E}\left[\left(\varepsilon_{i,j}^{(1)}\right)^2\right]\mathbb{E}\left[\left(1 + \frac{\eta_{i,j}}{p_{i,j}}\right)^2\right] \\ &= \mathbb{E}\left[\left(\varepsilon_{i,j}^{(1)}\right)^2\right]\left[1 + \mathbb{E}\left[\frac{\eta_{i,j}^2}{p_{i,j}^2}\right]\right] \end{aligned}$$

$$\stackrel{(a)}{=} (\sigma_{i,j}^{(1)})^2 \left[ 1 + \frac{p_{i,j}(1-p_{i,j})}{p_{i,j}^2} \right] = \frac{(\sigma_{i,j}^{(1)})^2}{p_{i,j}},$$

where (a) follows because  $\mathbb{E}[\eta_{i,j}^2] = \text{Var}(\eta_{i,j}) = p_{i,j}(1-p_{i,j})$  from Eq. (3), and  $\mathbb{E}[(\varepsilon_{i,j}^{(1)})^2] = \text{Var}(\varepsilon_{i,j}^{(1)}) = (\sigma_{i,j}^{(1)})^2$  from condition (C3). A similar argument establishes Eq. (A.25). Applying the same observations as above, we obtain Eq. (A.26) as follows,

$$\begin{aligned} \text{Cov}\left(\varepsilon_{i,j}^{(1)}\left(1 + \frac{\eta_{i,j}}{p_{i,j}}\right), \varepsilon_{i,j}^{(0)}\left(1 - \frac{\eta_{i,j}}{1-p_{i,j}}\right)\right) &= \mathbb{E}\left[\varepsilon_{i,j}^{(1)}\left(1 + \frac{\eta_{i,j}}{p_{i,j}}\right) \times \varepsilon_{i,j}^{(0)}\left(1 - \frac{\eta_{i,j}}{1-p_{i,j}}\right)\right] \\ &\stackrel{(a)}{=} \mathbb{E}\left[\left(1 + \frac{\eta_{i,j}}{p_{i,j}}\right)\left(1 - \frac{\eta_{i,j}}{1-p_{i,j}}\right)\right] \mathbb{E}[\varepsilon_{i,j}^{(1)}\varepsilon_{i,j}^{(0)}] \\ &= \left(1 - \mathbb{E}\left[\frac{\eta_{i,j}^2}{p_{i,j}(1-p_{i,j})}\right]\right) \mathbb{E}[\varepsilon_{i,j}^{(1)}\varepsilon_{i,j}^{(0)}] \\ &\stackrel{(b)}{=} 0 \cdot \mathbb{E}[\varepsilon_{i,j}^{(1)}\varepsilon_{i,j}^{(0)}] = 0, \end{aligned}$$

where (a) follows because  $(\varepsilon_{i,j}^{(0)}, \varepsilon_{i,j}^{(1)}) \perp\!\!\!\perp \eta_{i,j}$  from Assumption 2 and (b) follows because  $\mathbb{E}[\eta_{i,j}^2] = \text{Var}(\eta_{i,j}) = p_{i,j}(1-p_{i,j})$  from Eq. (3).

#### C.2.2. Proof of Lyapunov's condition with $\omega = 1$

We have

$$\begin{aligned} \frac{\sum_{i \in [N]} \mathbb{E}[|\mathbb{Z}_{i,j}^{\text{DR}}|^3]}{(\sum_{i \in [N]} \text{Var}(\mathbb{Z}_{i,j}^{\text{DR}}))^{3/2}} &= \frac{1}{N^{3/2}} \cdot \frac{\sum_{i \in [N]} \mathbb{E}[|\mathbb{Z}_{i,j}^{\text{DR}}|^3]}{(\frac{1}{N} \sum_{i \in [N]} \text{Var}(\mathbb{Z}_{i,j}^{\text{DR}}))^{3/2}} \stackrel{(a)}{=} \frac{1}{N^{3/2}} \cdot \frac{\sum_{i \in [N]} \mathbb{E}[|\mathbb{Z}_{i,j}^{\text{DR}}|^3]}{(\bar{\sigma}_j)^{3/2}} \\ &\stackrel{(b)}{\leq} \frac{1}{N^{3/2}} \cdot \frac{\sum_{i \in [N]} \mathbb{E}[|\mathbb{Z}_{i,j}^{\text{DR}}|^3]}{c_1^{3/2}} \\ &\stackrel{(c)}{\leq} \frac{1}{N^{1/2}} \cdot \frac{c_2}{c_1^{3/2}}, \end{aligned} \tag{A.27}$$

where (a) follows by putting together Eqs. (22) and (A.22), (b) follows because  $\bar{\sigma}_j \geq c_1 > 0$  as per condition (C3), (c) follows because the absolute third moments of subExponential random variables are bounded, after noting that  $\mathbb{Z}_{i,j}^{\text{DR}}$  is a subExponential random variable. Then, condition (A.21) holds for  $\omega = 1$  as the right hand side of Eq. (A.27) goes to 0 as  $N \rightarrow \infty$ .

## D. Proof of Proposition 1 (20): Finite Sample Guarantees for OI

Fix any  $j \in [M]$ . Recall the definitions Eqs. (5) and (9) of the parameters  $\text{ATE}_{\cdot,j}$  and corresponding outcome imputation estimates  $\widehat{\text{ATE}}_{\cdot,j}^{\text{OI}}$ . The error  $\Delta \text{ATE}_{\cdot,j}^{\text{OI}} = \widehat{\text{ATE}}_{\cdot,j}^{\text{OI}} - \text{ATE}_{\cdot,j}$  can be re-expressed as

$$\Delta \text{ATE}_{\cdot,j}^{\text{OI}} = \frac{1}{N} \sum_{i \in [N]} \left( \widehat{\theta}_{i,j}^{(1)} - \widehat{\theta}_{i,j}^{(0)} \right) - \frac{1}{N} \sum_{i \in [N]} \left( \theta_{i,j}^{(1)} - \theta_{i,j}^{(0)} \right)$$



$$= \frac{1}{N} \sum_{i \in [N]} \left( (\widehat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) - (\widehat{\theta}_{i,j}^{(0)} - \theta_{i,j}^{(0)}) \right).$$

Using the triangle inequality, we have

$$|\Delta \text{ATE}_{\cdot,j}^{\text{OI}}| \leq \frac{1}{N} \left| \sum_{i \in [N]} (\widehat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) \right| + \frac{1}{N} \left| \sum_{i \in [N]} (\widehat{\theta}_{i,j}^{(0)} - \theta_{i,j}^{(0)}) \right|. \quad (\text{A.28})$$

Consider any  $a \in \{0, 1\}$ . We claim that

$$\frac{1}{N} \left| \sum_{i \in [N]} (\widehat{\theta}_{i,j}^{(a)} - \theta_{i,j}^{(a)}) \right| \leq \mathcal{E}(\widehat{\Theta}^{(a)}). \quad (\text{A.29})$$

The proof is complete by putting together Eqs. (A.28) and (A.29).

**Proof of Eq. (A.29)** Fix any  $a \in \{0, 1\}$ . Using the Cauchy-Schwarz inequality, we have

$$\frac{1}{N} \left| \sum_{i \in [N]} (\widehat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)}) \right| \leq \frac{1}{N} \|\mathbf{1}_N\|_2 \|\widehat{\Theta}_{\cdot,j}^{(1)} - \Theta_{\cdot,j}^{(1)}\|_2 = \frac{1}{\sqrt{N}} \|\widehat{\Theta}_{\cdot,j}^{(1)} - \Theta_{\cdot,j}^{(1)}\|_2 \leq \frac{1}{\sqrt{N}} \|\widehat{\Theta}^{(1)} - \Theta^{(1)}\|_{1,2}.$$

The proof is complete by using the notation in Eq. (16).

## E. Proof of Proposition 1 (21): Finite Sample Guarantees for IPW

Fix any  $j \in [M]$ . Recall the definitions Eqs. (5) and (10) of the parameters  $\text{ATE}_{\cdot,j}$  and corresponding inverse probability weighting estimates  $\widehat{\text{ATE}}_{\cdot,j}^{\text{IPW}}$ . The error  $\Delta \text{ATE}_{\cdot,j}^{\text{IPW}} = \widehat{\text{ATE}}_{\cdot,j}^{\text{IPW}} - \text{ATE}_{\cdot,j}$  can be re-expressed as

$$\begin{aligned} \Delta \text{ATE}_{\cdot,j}^{\text{IPW}} &= \frac{1}{N} \sum_{i \in [N]} \left( \frac{y_{i,j} a_{i,j}}{\widehat{p}_{i,j}} - \frac{y_{i,j} (1 - a_{i,j})}{1 - \widehat{p}_{i,j}} \right) - \frac{1}{N} \sum_{i \in [N]} \left( \theta_{i,j}^{(1)} - \theta_{i,j}^{(0)} \right) \\ &= \frac{1}{N} \sum_{i \in [N]} \left( \left( \frac{y_{i,j} a_{i,j}}{\widehat{p}_{i,j}} - \theta_{i,j}^{(1)} \right) - \left( \frac{y_{i,j} (1 - a_{i,j})}{1 - \widehat{p}_{i,j}} - \theta_{i,j}^{(0)} \right) \right) \\ &\stackrel{(a)}{=} \frac{1}{N} \sum_{i \in [N]} \left( \mathbb{T}_{i,j}^{(1,\text{IPW})} + \mathbb{T}_{i,j}^{(0,\text{IPW})} \right), \end{aligned} \quad (\text{A.30})$$

where (a) follows after defining  $\mathbb{T}_{i,j}^{(1,\text{IPW})} \triangleq y_{i,j} a_{i,j} / \widehat{p}_{i,j} - \theta_{i,j}^{(1)}$  and  $\mathbb{T}_{i,j}^{(0,\text{IPW})} \triangleq \theta_{i,j}^{(0)} - y_{i,j} (1 - a_{i,j}) / (1 - \widehat{p}_{i,j})$  for every  $(i, j) \in [N] \times [M]$ . Then, we have

$$\begin{aligned} \mathbb{T}_{i,j}^{(1,\text{IPW})} &= \frac{y_{i,j} a_{i,j}}{\widehat{p}_{i,j}} - \theta_{i,j}^{(1)} \\ &\stackrel{(a)}{=} \frac{(\theta_{i,j}^{(1)} + \varepsilon_{i,j}^{(1)})(p_{i,j} + \eta_{i,j})}{\widehat{p}_{i,j}} - \theta_{i,j}^{(1)} \end{aligned}$$

$$\begin{aligned}
&= \theta_{i,j}^{(1)} \left( \frac{p_{i,j} + \eta_{i,j}}{\hat{p}_{i,j}} - 1 \right) + \varepsilon_{i,j}^{(1)} \left( \frac{p_{i,j} + \eta_{i,j}}{\hat{p}_{i,j}} \right) \\
&= \frac{\theta_{i,j}^{(1)} (p_{i,j} - \hat{p}_{i,j})}{\hat{p}_{i,j}} + \frac{\theta_{i,j}^{(1)} \eta_{i,j}}{\hat{p}_{i,j}} + \frac{\varepsilon_{i,j}^{(1)} p_{i,j}}{\hat{p}_{i,j}} + \frac{\varepsilon_{i,j}^{(1)} \eta_{i,j}}{\hat{p}_{i,j}}, \tag{A.31}
\end{aligned}$$

where (a) follows from Eqs. (1) to (3). A similar derivation for  $a = 0$  implies that

$$\begin{aligned}
\mathbb{T}_{i,j}^{(0, \text{IPW})} &= \theta_{i,j}^{(0)} - \frac{y_{i,j}(1 - a_{i,j})}{1 - \hat{p}_{i,j}} \\
&= -\frac{\theta_{i,j}^{(0)}(1 - p_{i,j} - (1 - \hat{p}_{i,j}))}{1 - \hat{p}_{i,j}} - \frac{\theta_{i,j}^{(0)}(-\eta_{i,j})}{1 - \hat{p}_{i,j}} - \frac{\varepsilon_{i,j}^{(0)}(1 - p_{i,j})}{1 - \hat{p}_{i,j}} - \frac{\varepsilon_{i,j}^{(0)}(-\eta_{i,j})}{1 - \hat{p}_{i,j}} \\
&= \frac{\theta_{i,j}^{(0)}(p_{i,j} - \hat{p}_{i,j})}{1 - \hat{p}_{i,j}} + \frac{\theta_{i,j}^{(0)}\eta_{i,j}}{1 - \hat{p}_{i,j}} - \frac{\varepsilon_{i,j}^{(0)}(1 - p_{i,j})}{1 - \hat{p}_{i,j}} + \frac{\varepsilon_{i,j}^{(0)}\eta_{i,j}}{1 - \hat{p}_{i,j}}.
\end{aligned}$$

Consider any  $a \in \{0, 1\}$  and any  $\delta \in (0, 1)$ . We claim that, with probability at least  $1 - 6\delta$ ,

$$\frac{1}{N} \left| \sum_{i \in [N]} \mathbb{T}_{i,j}^{(a, \text{IPW})} \right| \leq \frac{2}{\lambda} \|\Theta^{(a)}\|_{\max} \cdot \mathcal{E}(\hat{P}) + \frac{2\sqrt{c\ell_\delta}}{\lambda\sqrt{\ell_1 N}} \|\Theta^{(a)}\|_{\max} + \frac{2\bar{\sigma}\sqrt{c\ell_\delta}}{\lambda\sqrt{N}} + \frac{2\bar{\sigma}m(c\ell_\delta)}{\lambda\sqrt{\ell_1 N}}. \tag{A.32}$$

where recall that  $m(c\ell_\delta) = \max(c\ell_\delta, \sqrt{c\ell_\delta})$ . We provide a proof of this claim at the end of this section. Applying triangle inequality in Eq. (A.30) and using Eq. (A.32) with a union bound, we obtain that

$$|\Delta \text{ATE}_{i,j}^{\text{IPW}}| \leq \frac{2}{\lambda} \|\Theta\|_{\max} \cdot \mathcal{E}(\hat{P}) + \frac{2\sqrt{c\ell_\delta}}{\lambda\sqrt{\ell_1 N}} \|\Theta\|_{\max} + \frac{4\bar{\sigma}\sqrt{c\ell_\delta}}{\lambda\sqrt{N}} + \frac{4\bar{\sigma}m(c\ell_\delta)}{\lambda\sqrt{\ell_1 N}},$$

with probability at least  $1 - 12\delta$ . The claim in Eq. (21) follows by re-parameterizing  $\delta$ .

**Proof of Eq. (A.32).** This proof follows a very similar road map to that used for establishing the inequality in display (A.4). Recall the partitioning of the units  $[N]$  into  $\mathcal{R}_0$  and  $\mathcal{R}_1$  from Assumption 4. Condition on this partition. Now, to enable the application of concentration bounds, we split the summation over  $i \in [N]$  in the left hand side of Eq. (A.32) into two parts—one over  $i \in \mathcal{R}_0$  and the other over  $i \in \mathcal{R}_1$ —such that the noise terms are independent of the estimates of  $\Theta^{(0)}, \Theta^{(1)}, P$  in each of these parts as in Eqs. (14) and (15).

Fix  $a = 1$  and note that  $|\sum_{i \in [N]} \mathbb{T}_{i,j}^{(1, \text{IPW})}| \leq |\sum_{i \in \mathcal{R}_0} \mathbb{T}_{i,j}^{(1, \text{IPW})}| + |\sum_{i \in \mathcal{R}_1} \mathbb{T}_{i,j}^{(1, \text{IPW})}|$ . Fix any  $s \in \{0, 1\}$ . Then, Eq. (A.31) and triangle inequality imply that

$$\left| \sum_{i \in \mathcal{R}_s} \mathbb{T}_{i,j}^{(1, \text{IPW})} \right| \leq \left| \sum_{i \in \mathcal{R}_s} \frac{\theta_{i,j}^{(1)}(p_{i,j} - \hat{p}_{i,j})}{\hat{p}_{i,j}} \right| + \left| \sum_{i \in \mathcal{R}_s} \frac{\theta_{i,j}^{(1)}\eta_{i,j}}{\hat{p}_{i,j}} \right| + \left| \sum_{i \in \mathcal{R}_s} \frac{\varepsilon_{i,j}^{(1)}p_{i,j}}{\hat{p}_{i,j}} \right| + \left| \sum_{i \in \mathcal{R}_s} \frac{\varepsilon_{i,j}^{(1)}\eta_{i,j}}{\hat{p}_{i,j}} \right|. \tag{A.33}$$

Next, note that the decomposition in Eq. (A.33) is identical to the one in Eq. (A.6), except for the fact when compared to Eq. (A.6), the first two terms in Eq. (A.33) have a

factor of  $\theta_{i,j}^{(1)}$  instead of  $(\hat{\theta}_{i,j}^{(1)} - \theta_{i,j}^{(1)})$ . As a result, mimicking steps used to derive Eq. (A.12), we obtain the following bound, with probability at least  $1 - 3\delta$ ,

$$\begin{aligned} \frac{1}{N} \left| \sum_{i \in \mathcal{R}_s} \mathbb{T}_{i,j}^{(1, \text{IPW})} \right| &\leq \frac{1}{\bar{\lambda}N} \|\Theta^{(1)}\|_{1,2} \|\hat{P} - P\|_{1,2} + \frac{\sqrt{cl_\delta}}{\bar{\lambda}\sqrt{\ell_1}N} \|\Theta^{(1)}\|_{1,2} + \frac{\bar{\sigma}\sqrt{cl_\delta}}{\bar{\lambda}N} \|P\|_{1,2} + \frac{\bar{\sigma}m(cl_\delta)}{\bar{\lambda}\sqrt{\ell_1}N} \|\mathbf{1}\|_{1,2} \\ &\stackrel{(a)}{\leq} \frac{1}{\bar{\lambda}\sqrt{N}} \|\Theta^{(1)}\|_{\max} \|\hat{P} - P\|_{1,2} + \frac{\sqrt{cl_\delta}}{\bar{\lambda}\sqrt{\ell_1}N} \|\Theta^{(1)}\|_{\max} + \frac{\bar{\sigma}\sqrt{cl_\delta}}{\bar{\lambda}\sqrt{N}} + \frac{\bar{\sigma}m(cl_\delta)}{\bar{\lambda}\sqrt{\ell_1}N}, \\ &\stackrel{(b)}{\leq} \frac{1}{\bar{\lambda}} \|\Theta^{(1)}\|_{\max} \cdot \mathcal{E}(\hat{P}) + \frac{\sqrt{cl_\delta}}{\bar{\lambda}\sqrt{\ell_1}N} \|\Theta^{(1)}\|_{\max} + \frac{\bar{\sigma}\sqrt{cl_\delta}}{\bar{\lambda}\sqrt{N}} + \frac{\bar{\sigma}m(cl_\delta)}{\bar{\lambda}\sqrt{\ell_1}N}, \end{aligned} \quad (\text{A.34})$$

where (a) follows because  $\|\Theta^{(1)}\|_{1,2} \leq \sqrt{N} \|\Theta^{(1)}\|_{\max}$ ,  $\|P\|_{1,2} \leq \sqrt{N}$  and  $\|\mathbf{1}\|_{1,2} = \sqrt{N}$ , and (b) follows from Eq. (16). Then, the claim in Eq. (A.32) follows for  $a = 1$  by using Eq. (A.34) and applying a union bound over  $s \in \{0, 1\}$ . The proof of Eq. (A.32) for  $a = 0$  follows similarly.

## F. Proofs of Propositions 2 and 3

In Appendix F.1, we prove Proposition 2, i.e., we show that the estimates of  $P$ ,  $\Theta^{(0)}$ , and  $\Theta^{(1)}$  generated by **Cross-Fitted-MC** satisfy Assumption 4. Next, we prove Proposition 3 implying that the estimates of  $P$ ,  $\Theta^{(0)}$ , and  $\Theta^{(1)}$  generated by **Cross-Fitted-SVD** satisfy the condition (C2) in Theorem 2 as long as  $\sqrt{N}/M = o(1)$ .

### F.1. Proof of Proposition 2: Guarantees for Cross-Fitted-MC

Consider any matrix completion algorithm MC and any block partition  $\mathcal{P}$  of the set  $[N] \times [M]$  into four blocks as in Assumption 4. Fix any  $\mathcal{I} \in \mathcal{P}$ .

Consider  $\hat{P}$  in Eq. (27). To see why  $\hat{P}_{\mathcal{I}} \perp\!\!\!\perp W_{\mathcal{I}}$ , note that  $A \otimes \mathbf{1}^{-\mathcal{I}}$  depends only on  $W \setminus W_{\mathcal{I}}$  which is independent of  $W_{\mathcal{I}}$  from Assumption 2(c). Further,  $\hat{P}_{\mathcal{I}} \perp\!\!\!\perp E_{\mathcal{I}}^{(0)}, E_{\mathcal{I}}^{(1)}$  holds since  $W \setminus W_{\mathcal{I}}$  is independent of  $(E_{\mathcal{I}}^{(0)}, E_{\mathcal{I}}^{(1)})$  from Assumption 2(b). Overall, we conclude  $\hat{P}_{\mathcal{I}} \perp\!\!\!\perp (W_{\mathcal{I}}, E_{\mathcal{I}}^{(0)}, E_{\mathcal{I}}^{(1)})$ .

Next, consider  $\hat{\Theta}^{(0)}$  and  $\hat{\Theta}^{(1)}$  defined in Eqs. (25) and (26), respectively. Fix any  $a \in \{0, 1\}$ . To see why  $\hat{\Theta}_{\mathcal{I}}^{(a)} \perp\!\!\!\perp W_{\mathcal{I}}$ , note that  $\hat{\Theta}_{\mathcal{I}}^{(a)}$  depends on  $Y^{(a), \text{obs}} \otimes \mathbf{1}^{-\mathcal{I}}$ , which in turn depends on (i)  $W \setminus W_{\mathcal{I}}$  and (ii)  $E^{(a)} \setminus E_{\mathcal{I}}^{(a)}$ , each of which are independent of  $W_{\mathcal{I}}$  due to Assumption 2(c) and Assumption 2(b), respectively.

### F.2. Proof of Proposition 3: Guarantees for Cross-Fitted-SVD

To prove this result, we first derive a corollary of Lemma A.1 in Bai and Ng (2021) for a generic matrix of interest  $T$ , such that  $S = (T + H) \otimes F$ , and apply it to  $P$ ,  $\Theta^{(0)} \odot (\mathbf{1} - P)$ , and  $\Theta^{(1)} \odot P$ . We impose the following restrictions on  $T$  and  $H$ .

**Assumption 8.** *There exist a constant  $r_T \in [\min\{N, M\}]$  and a collection of latent factors*

$$\tilde{U} \in \mathbb{R}^{N \times r_T} \quad \text{and} \quad \tilde{V} \in \mathbb{R}^{M \times r_T},$$

*such that,*

- (a)  $T$  satisfies the factorization:  $T = \tilde{U}\tilde{V}^\top$ ,
- (b)  $\|\tilde{U}\|_{2,\infty} \leq c$  and  $\|\tilde{V}\|_{2,\infty} \leq c$  for some positive constant  $c$ , and
- (c)  $N^{-1}\tilde{U}^\top\tilde{U}$  and  $M^{-1}\tilde{V}^\top\tilde{V}$  are positive definite matrices.

**Assumption 9.** The noise matrix  $H$  is such that,

- (a)  $\{h_{i,j} : i \in [N], j \in [M]\}$  are zero-mean subExponential with the subExponential norm bounded by a constant  $\bar{\sigma}$ ,
- (b)  $\sum_{j' \in [M]} |\mathbb{E}[h_{i,j}h_{i,j'}]| \leq c$  for every  $i \in [N]$  and  $j \in [M]$ , and
- (c)  $\{H_{i,\cdot} : i \in [N]\}$  are mutually independent (across  $i$ ).

The next result characterizes the entry-wise error in recovering the missing entries of a matrix where all entries in one block are deterministically missing (see the discussion in Section 5.1) using the TW algorithm (summarized in Section 5.2.1). Its proof, essentially established as a corollary of Bai and Ng (2021, Lemma A.1), is provided in Appendix F.3.

**Corollary F.1.** Consider a matrix of interest  $T$  that satisfies Assumption 8 and a noise matrix  $H$  that satisfies Assumption 9. Let  $S \in \{\mathbb{R}, ?\}^{N \times M}$  be the observed matrix as in Eq. (6). Let  $\mathcal{R}_{\text{obs}} \subseteq [N]$  and  $\mathcal{C}_{\text{obs}} \subseteq [M]$  denote the set of rows and columns of  $S$ , respectively, with all entries observed. Suppose the mask matrix  $F$  is such that each  $i \in [N]$  belongs to  $\mathcal{R}_{\text{obs}}$  with probability  $1/2$  and each  $j \in [M]$  belongs to  $\mathcal{C}_{\text{obs}}$  with probability  $1/2$ . Let  $\mathcal{I} = \mathcal{R}_{\text{miss}} \times \mathcal{C}_{\text{miss}}$  where  $\mathcal{R}_{\text{miss}} \triangleq [N] \setminus \mathcal{R}_{\text{obs}}$  and  $\mathcal{C}_{\text{miss}} \triangleq [M] \setminus \mathcal{C}_{\text{obs}}$ . Then,  $\text{TW}_{r_T}$  produces an estimate  $\hat{T}_{\mathcal{I}}$  of  $T_{\mathcal{I}}$  such that

$$\|\hat{T}_{\mathcal{I}} - T_{\mathcal{I}}\|_{\max} = O_p\left(\frac{1}{\sqrt{N}} + \frac{1}{\sqrt{M}}\right),$$

as  $N, M \rightarrow \infty$ .

Given this corollary, we now complete the proof of Proposition 3. Consider the partition  $\mathcal{P}$  in step 2 of **Cross-Fitted-SVD** and fix any  $\mathcal{I} \in \mathcal{P}$ . Recall that **Cross-Fitted-SVD** applies TW on  $P \otimes \mathbf{1}^{-\mathcal{I}}$ ,  $Y^{(0),\text{full}} \otimes \mathbf{1}^{-\mathcal{I}}$ , and  $Y^{(1),\text{full}} \otimes \mathbf{1}^{-\mathcal{I}}$ , and note that  $\mathbf{1}^{-\mathcal{I}}$  satisfies the requirement on the mask matrix in Corollary F.1.

#### F.2.1. Estimating $P$ .

Consider estimating  $P$  using **Cross-Fitted-SVD**. To apply Corollary F.1, we use Assumptions 5 and 6 to note that  $P$  satisfies Assumption 8 with rank parameter  $r_p$ . Then, we use Eq. (3) and Assumption 2(b) to note that  $W$  satisfies Assumption 9. Step 3 of **Cross-Fitted-SVD** can be rewritten as  $\hat{P} = \text{Proj}_{\hat{\lambda}}(\bar{P})$  and  $\bar{P} = \text{Cross-Fitted-MC}(\text{TW}_{r_1}, A, \mathcal{P})$  where  $r_1 = r_p$ . Then,

$$\|\hat{P}_{\mathcal{I}} - P_{\mathcal{I}}\|_{\max} \stackrel{(a)}{\leq} \|\bar{P}_{\mathcal{I}} - P_{\mathcal{I}}\|_{\max} \stackrel{(b)}{=} O_p\left(\frac{1}{\sqrt{N}} + \frac{1}{\sqrt{M}}\right),$$

where (a) follows from Assumptions 1 and 3, and the definition of  $\text{Proj}_{\bar{\lambda}}(\cdot)$ , and (b) follows from Corollary F.1. Applying a union bound over all  $\mathcal{I} \in \mathcal{P}$ , we have

$$\mathcal{E}(\hat{P}) \stackrel{(a)}{\leq} \|\hat{P} - P\|_{\max} = O_p\left(\frac{1}{\sqrt{N}} + \frac{1}{\sqrt{M}}\right), \quad (\text{A.35})$$

where (a) follows from the definition of  $L_{1,2}$  norm.

*F.2.2. Estimating  $\Theta^{(0)}$  and  $\Theta^{(1)}$ .*

For every  $a \in \{0, 1\}$ , we show that

$$\mathcal{E}(\hat{\Theta}^{(a)}) = O_p\left(\frac{1}{\sqrt{N}} + \frac{1}{\sqrt{M}}\right). \quad (\text{A.36})$$

We focus on  $a = 1$  noting that the proof for  $a = 0$  is analogous. We split the proof in two cases: (i)  $\|(\hat{\Theta}^{(1)} - \Theta^{(1)}) \odot \hat{P}\|_{\max} \leq \|\Theta^{(1)} \odot (\hat{P} - P)\|_{\max}$  and (ii)  $\|(\hat{\Theta}^{(1)} - \Theta^{(1)}) \odot \hat{P}\|_{\max} \geq \|\Theta^{(1)} \odot (\hat{P} - P)\|_{\max}$ .

In the first case, we have

$$\begin{aligned} \bar{\lambda} \|\hat{\Theta}^{(1)} - \Theta^{(1)}\|_{\max} &\stackrel{(a)}{\leq} \|(\hat{\Theta}^{(1)} - \Theta^{(1)}) \odot \hat{P}\|_{\max} \leq \|\Theta^{(1)} \odot (\hat{P} - P)\|_{\max} \\ &\stackrel{(b)}{\leq} \|\Theta^{(1)}\|_{\max} \|\hat{P} - P\|_{\max}, \end{aligned} \quad (\text{A.37})$$

where (a) follows from Assumption 3 and (b) follows from the definition of  $\|\Theta^{(1)}\|_{\max}$ . Then,

$$\mathcal{E}(\hat{\Theta}^{(1)}) \stackrel{(a)}{\leq} \|\hat{\Theta}^{(1)} - \Theta^{(1)}\|_{\max} \stackrel{(b)}{\leq} \frac{\|\Theta^{(1)}\|_{\max}}{\bar{\lambda}} \|\hat{P} - P\|_{\max} \stackrel{(c)}{=} \frac{\|\Theta^{(1)}\|_{\max}}{\bar{\lambda}} O_p\left(\frac{1}{\sqrt{N}} + \frac{1}{\sqrt{M}}\right),$$

where (a) follows from the definition of  $L_{1,2}$  norm, (b) follows from Eq. (A.37), and (c) follows from Eq. (A.35). Then, Eq. (A.36) follows as  $1/\bar{\lambda}$  and  $\|\Theta^{(1)}\|_{\max}$  are assumed to be bounded.

In the second case, using Eqs. (2) and (3) to expand  $Y^{(1),\text{full}}$ , we have

$$Y^{(1),\text{full}} = \Theta^{(1)} \odot P + \Theta^{(1)} \odot \eta + \varepsilon^{(1)} \odot P + \varepsilon^{(1)} \odot \eta.$$

Next, we utilize two claims proven in Appendices F.2.3 and F.2.4 respectively:  $\Theta^{(1)} \odot P$  satisfies Assumption 8 with rank parameter  $r_{\theta_1} r_p$  and

$$\bar{\varepsilon}^{(1)} \triangleq \Theta^{(1)} \odot \eta + \varepsilon^{(1)} \odot P + \varepsilon^{(1)} \odot \eta,$$

satisfies Assumption 9.

Now, note that step 6 of **Cross-Fitted-SVD** can be rewritten as  $\hat{\Theta}^{(1)} = \bar{\Theta}^{(1)} \odot \hat{P}$  and  $\bar{\Theta}^{(1)} = \text{Cross-Fitted-MC}(\text{TW}_{r_3}, Y^{(1),\text{full}}, \mathcal{P})$  where  $r_3 = r_{\theta_1} r_p$ . Then, from Corollary F.1,

$$\|\bar{\Theta}^{(1)} - \Theta^{(1)} \odot P\|_{\max} = O_p\left(\frac{1}{\sqrt{N}} + \frac{1}{\sqrt{M}}\right).$$

Applying a union bound over all  $\mathcal{I} \in \mathcal{P}$  and noting that  $\bar{\Theta}^{(1)} = \hat{\Theta}^{(1)} \odot \hat{P}$ , we have

$$\|\hat{\Theta}^{(1)} \odot \hat{P} - \Theta^{(1)} \odot P\|_{\max} = O_p\left(\frac{1}{\sqrt{N}} + \frac{1}{\sqrt{M}}\right). \quad (\text{A.38})$$

The left hand side of Eq. (A.38) can be written as,

$$\begin{aligned} \|\hat{\Theta}^{(1)} \odot \hat{P} - \Theta^{(1)} \odot P\|_{\max} &= \|\hat{\Theta}^{(1)} \odot \hat{P} - \Theta^{(1)} \odot \hat{P} + \Theta^{(1)} \odot \hat{P} - \Theta^{(1)} \odot P\|_{\max} \\ &\stackrel{(a)}{\geq} \|(\hat{\Theta}^{(1)} - \Theta^{(1)}) \odot \hat{P}\|_{\max} - \|\Theta^{(1)} \odot (\hat{P} - P)\|_{\max} \\ &\stackrel{(b)}{\geq} \bar{\lambda} \|\hat{\Theta}^{(1)} - \Theta^{(1)}\|_{\max} - \|\Theta^{(1)}\|_{\max} \|\hat{P} - P\|_{\max}, \end{aligned} \quad (\text{A.39})$$

where (a) follows from triangle inequality as  $\|(\hat{\Theta}^{(1)} - \Theta^{(1)}) \odot \hat{P}\|_{\max} \geq \|\Theta^{(1)} \odot (\hat{P} - P)\|_{\max}$  and (b) follows from Assumption 3 and the definition of  $\|\Theta^{(1)}\|_{\max}$ . Then,

$$\begin{aligned} \mathcal{E}(\hat{\Theta}^{(1)}) &\stackrel{(a)}{\leq} \|\hat{\Theta}^{(1)} - \Theta^{(1)}\|_{\max} \stackrel{(b)}{\leq} \frac{1}{\bar{\lambda}} \|\hat{\Theta}^{(1)} \odot \hat{P} - \Theta^{(1)} \odot P\|_{\max} + \frac{\|\Theta^{(1)}\|_{\max}}{\bar{\lambda}} \|\hat{P} - P\|_{\max} \\ &\stackrel{(b)}{=} \frac{1}{\bar{\lambda}} O_p\left(\frac{1}{\sqrt{N}} + \frac{1}{\sqrt{M}}\right) + \frac{\|\Theta^{(1)}\|_{\max}}{\bar{\lambda}} O_p\left(\frac{1}{\sqrt{N}} + \frac{1}{\sqrt{M}}\right), \end{aligned}$$

where (a) follows from the definition of  $L_{1,2}$  norm, (b) follows from Eq. (A.39), and (c) follows from Eqs. (A.35) and (A.38). Then, Eq. (A.36) follows as  $1/\bar{\lambda}$  and  $\|\Theta^{(1)}\|_{\max}$  are assumed to be bounded.

*F.2.3. Proof that  $\Theta^{(0)} \odot (\mathbf{1} - P)$  and  $\Theta^{(1)} \odot P$  satisfy Assumption 8.*

Recall the factors  $\bar{U}^{(0)}$  and  $\bar{V}^{(0)}$  of  $\Theta^{(0)} \odot (\mathbf{1} - P)$ , and  $\bar{U}^{(1)}$  and  $\bar{V}^{(1)}$  of  $\Theta^{(1)} \odot P$  from Section 5.3. Then, Assumption 8(a) holds from Eq. (29). Next, we note that

$$\|\bar{U}^{(1)}\|_{2,\infty} = \|U * U^{(1)}\|_{2,\infty} \stackrel{(a)}{=} \max_{i \in [N]} \sqrt{\sum_{j \in [r_p]} u_{i,j}^2 \sum_{j' \in [r_{\theta_1}]} (u_{i,j'}^{(1)})^2} \leq \|U\|_{2,\infty} \|U^{(1)}\|_{2,\infty} \stackrel{(b)}{\leq} c,$$

where (a) follows from the definition of Khatri-Rao product (see Section 1), and (b) follows from Assumption 6. Then,  $\Theta^{(1)} \odot P$  satisfies Assumption 8(b) by using similar arguments on  $\bar{V}^{(1)}$ . Further,  $\Theta^{(0)} \odot (\mathbf{1} - P)$  satisfies Assumption 8(b) by noting that  $\|\bar{U}\|_{2,\infty}$  and  $\|\bar{V}\|_{2,\infty}$  are bounded whenever  $\|U\|_{2,\infty}$  and  $\|V\|_{2,\infty}$  are bounded, respectively. Finally, Assumption 8(c) holds from Assumption 6.

*F.2.4. Proof that  $\bar{\varepsilon}^{(1)}$  satisfies Assumption 9*

Recall that  $\bar{\varepsilon}^{(1)} \triangleq \Theta^{(1)} \odot \eta + \varepsilon^{(1)} \odot P + \varepsilon^{(1)} \odot \eta$ . Then, Assumption 9(a) holds as  $\bar{\varepsilon}^{(1)}$  is zero-mean from Assumption 2 and Eq. (3), and  $\bar{\varepsilon}^{(1)}$  is subExponential because  $\varepsilon_{i,j}^{(1)} \eta_{i,j}$  is a subExponential random variable Lemma A.3, every subGaussian random variable is subExponential random

variable, and sum of subExponential random variables is a subExponential random variable. Next, Assumption 9(b) holds as

$$\sum_{j' \in [M]} |\mathbb{E}[\bar{\varepsilon}_{i,j}^{(1)} \bar{\varepsilon}_{i,j'}^{(1)}]| \stackrel{(a)}{=} \sum_{j' \in [M]} |\mathbb{E}[\theta_{i,j}^{(1)} \theta_{i,j'}^{(1)} \eta_{i,j} \eta_{i,j'} + p_{i,j} p_{i,j'} \bar{\varepsilon}_{i,j}^{(1)} \bar{\varepsilon}_{i,j'}^{(1)} + \bar{\varepsilon}_{i,j}^{(1)} \bar{\varepsilon}_{i,j'}^{(1)} \eta_{i,j} \eta_{i,j'}]| \stackrel{(b)}{\leq} c,$$

where (a) follows from Assumption 2, and (b) follows from Assumptions 1, 2, and 7, Eq. (3), Lemma A.3, and because  $\|\Theta^{(1)}\|_{\max}$  are bounded. Finally, Assumption 9(b) holds from Assumptions 2 and 7.

### F.3. Proof of Corollary F.1

Corollary F.1 is a direct application of Bai and Ng (2021, Lemma A.1), specialized to our setting. Notably, Bai and Ng (2021) make four assumptions numbered A, B, C and D in their paper to establish the corresponding result. It remains to establish that the conditions assumed in Corollary F.1 imply the necessary conditions used in the proof of Bai and Ng (2021, Lemma A.1). First, note that due to the specific sampling assumed in defining the mask matrix in Corollary F.1, Bai and Ng (2021, Assumption D) holds immediately and Bai and Ng (2021, Assumption B) holds with high probability by Hoeffding's inequality.

It remains to show how Assumptions 8 and 9 imply the remainder of their assumptions, namely Bai and Ng (2021, Assumptions A and C). Before doing that, note that certain assumptions in Bai and Ng (2021) are not actually used in their proof of Lemma A.1 (or in the proof of other results used in that proof), namely, the distinct eigenvalue condition in Assumption A(a)(iii), the asymptotic normality conditions in Assumption A(c) and the asymptotic normality conditions in Assumption C. For completeness, the remaining relevant conditions from Bai and Ng (2021) are collected in the following two assumptions.

**Assumption 10** (Strong block factors). *Consider the latent factors  $\tilde{U} \in \mathbb{R}^{N \times r_T}$  and  $\tilde{V} \in \mathbb{R}^{M \times r_T}$  from Assumption 8. Define the following matrices:*

$$\tilde{U}^{\text{obs}} \triangleq \tilde{U}_{\mathcal{R}_{\text{obs}} \times [r_T]}, \quad \tilde{U}^{\text{miss}} \triangleq \tilde{U}_{\mathcal{R}_{\text{miss}} \times [r_T]}, \quad \tilde{V}^{\text{obs}} \triangleq \tilde{V}_{\mathcal{C}_{\text{obs}} \times [r_T]}, \quad \text{and} \quad \tilde{V}^{\text{miss}} \triangleq \tilde{V}_{\mathcal{C}_{\text{miss}} \times [r_T]},$$

where  $\mathcal{R}_{\text{obs}} \subseteq [N]$  and  $\mathcal{C}_{\text{obs}} \subseteq [M]$  denote the set of rows and columns of  $S$ , respectively, with all entries observed, and  $\mathcal{R}_{\text{miss}} \triangleq [N] \setminus \mathcal{R}_{\text{obs}}$  and  $\mathcal{C}_{\text{miss}} \triangleq [M] \setminus \mathcal{C}_{\text{obs}}$ . Then, the matrices defined below are positive definite:

$$\Sigma^{\tilde{U}, \text{obs}} \triangleq \frac{\tilde{U}^{\text{obs} \top} \tilde{U}^{\text{obs}}}{|\mathcal{R}_{\text{obs}}|}, \quad \Sigma^{\tilde{U}, \text{miss}} \triangleq \frac{\tilde{U}^{\text{miss} \top} \tilde{U}^{\text{miss}}}{|\mathcal{R}_{\text{miss}}|}, \quad \Sigma^{\tilde{V}, \text{obs}} \triangleq \frac{\tilde{V}^{\text{obs} \top} \tilde{V}^{\text{obs}}}{|\mathcal{C}_{\text{obs}}|}, \quad \text{and} \quad \Sigma^{\tilde{V}, \text{miss}} \triangleq \frac{\tilde{V}^{\text{miss} \top} \tilde{V}^{\text{miss}}}{|\mathcal{C}_{\text{miss}}|}.$$

**Assumption 11.** *The noise matrix  $H$  is such that,*

- (a)  $\max_{j \in [M]} \frac{1}{N} \sum_{j' \in [M]} |\sum_{i \in [N]} \mathbb{E}[h_{i,j} h_{i,j'}]| \leq c,$
- (b)  $\max_{j \in [M]} |\mathbb{E}[h_{i,j} h_{i',j}]| \leq c_{i,i'}$  and  $\max_{i \in [N]} \sum_{i' \in [N]} c_{i,i'} \leq c,$
- (c)  $\frac{1}{NM} \sum_{i,i' \in [N]} \sum_{j,j' \in [M]} |\mathbb{E}[h_{i,j} h_{i',j'}]| \leq c,$  and

$$(d) \max_{j,j' \in [M]} \frac{1}{N^2} \mathbb{E} \left[ \left| \sum_{i \in [N]} (h_{i,j} h_{i,j'} - \mathbb{E}[h_{i,j} h_{i,j'}]) \right|^4 \right].$$

Assumption 10 is a restatement of Bai and Ng (2021, Assumption C) (without the central limit theorems, which are not used in Bai and Ng (2021, Proof of Lemma A.1) as noted above). This condition ensures a strong factor structure on the sub-matrix corresponding to observed elements of  $S$  as well as on the sub-matrix corresponding to missing elements of  $S$ .

Assumption 11 is a restatement of the subset of conditions from Bai and Ng (2021, Assumption A) necessary in Bai and Ng (2021, proof of Lemma A.1) and it essentially requires weak dependence in the noise across measurements and across units. In particular, Assumption 11(a), (b), (c), and (d) correspond to Assumption A(b)(ii), (iii), (iv), (v), respectively, of Bai and Ng (2021). For the other conditions in Bai and Ng (2021, Assumption A), note that Assumption 8 above is equivalent to their Assumption A(a)(i) and (ii) of Bai and Ng (2021) when the factors are non-random as in this work. Similarly, Assumption 9(a) above is analogous to Assumption A(b)(i) of Bai and Ng (2021). Assumption A(b)(vi) of Bai and Ng (2021) is implied by their other Assumptions for non-random factors as stated in Bai (2003).

To establish Corollary F.1, it remains to establish that Assumptions 10 and 11 hold, which is done in Appendices F.3.1 and F.3.2 respectively.

### F.3.1. Assumption 10 holds

We show that  $\Sigma^{\tilde{U}, \text{obs}}$  is positive definite. The proof for  $\Sigma^{\tilde{U}, \text{miss}}$ ,  $\Sigma^{\tilde{U}, \text{obs}}$ , and  $\Sigma^{\tilde{U}, \text{miss}}$  being positive definite follows similarly. Define  $\Sigma^{\tilde{U}} \triangleq N^{-1} \tilde{U}^\top \tilde{U} \in \mathbb{R}^{r_T \times r_T}$ . From Weyl's inequality (Bhatia, 2007, Theorem. 8.2), we have the following for some  $c > 0$ :

$$\begin{aligned} \lambda_{\min}(\Sigma^{\tilde{U}, \text{obs}}) &\geq \lambda_{\min}(\Sigma^{\tilde{U}}) - \lambda_{\max}(\Sigma^{\tilde{U}} - \Sigma^{\tilde{U}, \text{obs}}) \stackrel{(a)}{\geq} c - \lambda_{\max}(\Sigma^{\tilde{U}} - \Sigma^{\tilde{U}, \text{obs}}) \\ &\geq c - |\lambda_{\max}(\Sigma^{\tilde{U}} - \Sigma^{\tilde{U}, \text{obs}})|, \end{aligned}$$

where (a) follows from Assumption 8(c) as  $\Sigma^{\tilde{U}}$  is positive definite. Now, it suffices to show that  $|\lambda_{\max}(\Sigma^{\tilde{U}} - \Sigma^{\tilde{U}, \text{obs}})| = o_p(1)$ .

Recall that the mask matrix  $F$  is such that each  $i \in [N]$  belongs to  $\mathcal{R}_{\text{obs}}$  with probability  $1/2$ . For every  $i \in [N]$ , let  $\mathbf{1}_i$  be an indicator random variable such that  $\mathbf{1}_i = 1$  if  $i \in \mathcal{R}_{\text{obs}}$  and  $\mathbf{1}_i = 0$  if  $i \notin \mathcal{R}_{\text{obs}}$ . Then, we express  $\Sigma^{\tilde{U}, \text{obs}}$  as follows,

$$\Sigma^{\tilde{U}, \text{obs}} = \frac{\sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top}{\sum_{i \in [N]} \mathbf{1}_i}. \quad (\text{A.40})$$

Then, we have

$$\begin{aligned} |\lambda_{\max}(\Sigma^{\tilde{U}} - \Sigma^{\tilde{U}, \text{obs}})| &\stackrel{(a)}{=} \left| \lambda_{\max} \left( \frac{\sum_{i \in [N]} \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top}{N} - \frac{\sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top}{\sum_{i \in [N]} \mathbf{1}_i} \right) \right| \\ &\stackrel{(b)}{\leq} \left| \lambda_{\max} \left( \frac{\sum_{i \in [N]} \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top}{N} - \frac{\sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top}{N/2} \right) \right| \end{aligned}$$



$$+ \left| \lambda_{\max} \left( \frac{\sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top}{N/2} - \frac{\sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top}{\sum_{i \in [N]} \mathbf{1}_i} \right) \right|, \quad (\text{A.41})$$

where (a) follows from Eq. (A.40) and the definition of  $\Sigma^{\tilde{U}}$ , and (b) follows from the triangle inequality on the operator norm after noting that the maximum eigenvalue of any symmetric matrix coincides with its operator norm. Next, we show that each term in Eq. (A.41) is  $o_p(1)$ .

**Proof that first term in Eq. (A.41) is  $o_p(1)$ .** We have

$$\lambda_{\max} \left( \frac{\sum_{i \in [N]} \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top}{N} - \frac{\sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top}{N/2} \right) = \lambda_{\max} \left( \frac{1}{N} \sum_{i \in [N]} (1 - 2\mathbf{1}_i) \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top \right). \quad (\text{A.42})$$

To bound Eq. (A.42), we apply Corollary A.3 with

$$X_i = \frac{1}{N} (1 - 2\mathbf{1}_i) \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top \quad \text{and} \quad A_i = \frac{1}{N} \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top. \quad (\text{A.43})$$

We note that, for every  $i \in [N]$ ,  $\mathbb{E}[X_i] = 0$  as  $\mathbb{E}[\mathbf{1}_i] = 1/2$  and  $A_i^2 - X_i^2$  is positive semi-definite as

$$X_i^2 = \frac{1}{N^2} (1 - 2\mathbf{1}_i)^2 \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top \stackrel{(a)}{=} \frac{1}{N^2} \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top = A_i^2,$$

where (a) follows because  $\mathbf{1}_i \in \{0, 1\}$ . We claim that  $|\lambda_{\max}(\sum_{i \in [n]} A_i^2)| \leq c^2/N$  for some  $c > 0$ . Then, using Corollary A.3, Eq. (A.42) is bounded as follows with probability at least  $1 - \delta$ ,

$$\lambda_{\max} \left( \frac{1}{N} \sum_{i \in [N]} (1 - 2\mathbf{1}_i) \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top \right) \leq \frac{\sqrt{c\ell_{2\delta/r_T}}}{\sqrt{N}}.$$

Therefore, the first term in Eq. (A.41) is  $o_p(1)$ . It remains to bound  $\lambda_{\max}(\sum_{i \in [N]} A_i^2)$ . We have

$$\begin{aligned} \left| \lambda_{\max} \left( \sum_{i \in [N]} A_i^2 \right) \right| &\stackrel{(a)}{=} \left| \max_{x \in \mathbb{R}^N: \|x\|_2=1} x^\top \left( \sum_{i \in [N]} A_i^2 \right) x \right| \\ &\stackrel{(b)}{=} \left| \max_{x \in \mathbb{R}^N: \|x\|_2=1} \frac{1}{N^2} \sum_{i \in [N]} x^\top \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top x \right| \\ &\stackrel{(c)}{=} \left| \max_{x \in \mathbb{R}^N: \|x\|_2=1} \frac{1}{N^2} \sum_{i \in [N]} \|\tilde{U}_{i,\cdot}\|_2 \cdot x^\top \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top x \right| \\ &\stackrel{(d)}{\leq} \max_{x \in \mathbb{R}^N: \|x\|_2=1} \frac{c}{N^2} \sum_{i \in [N]} |x^\top \tilde{U}_{i,\cdot}|^2 \stackrel{(e)}{\leq} \max_{x \in \mathbb{R}^N: \|x\|_2=1} \frac{c}{N^2} \sum_{i \in [N]} \|x\|_2 \|\tilde{U}_{i,\cdot}\|_2 \stackrel{(f)}{\leq} \frac{c^2}{N}, \end{aligned}$$

where (a) follows from the definition of the maximum eigenvalue of a matrix, (b) follows from Eq. (A.43), (c) follows because  $\tilde{U}_{i,\cdot} \in \mathbb{R}^{r_T \times 1}$ , (d) and (f) follow from Assumption 8(b), (e) follows from Cauchy-Schwarz inequality.

**Proof that second term in Eq. (A.41) is  $o_p(1)$ .** We have

$$\left| \lambda_{\max} \left( \frac{\sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top}{N/2} - \frac{\sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top}{\sum_{i \in [N]} \mathbf{1}_i} \right) \right| = \left| \lambda_{\max} \left( \frac{2}{N} \sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top \right) \right| \left| \frac{\sum_{i \in [N]} \mathbf{1}_i - N/2}{\sum_{i \in [N]} \mathbf{1}_i} \right|. \quad (\text{A.44})$$

To bound Eq. (A.44), we claim  $|\lambda_{\max}(\frac{1}{N} \sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top)| \leq c$  for some  $c > 0$ , and apply Corollary A.1 on the vector  $(\mathbf{1}_1 - 1/2, \dots, \mathbf{1}_N - 1/2)$ . We note that, for every  $i \in [N]$ ,  $\mathbf{1}_i - 1/2$  is zero-mean and subGaussian( $1/\sqrt{\ell_1}$ ) (see Example 2.5.8 in Vershynin (2018)). Then, with probability at least  $1 - \delta$ ,

$$\left| \sum_{i \in [N]} \mathbf{1}_i - N/2 \right| \leq \sqrt{\frac{c\ell_\delta N}{\ell_1}}. \quad (\text{A.45})$$

Using Eq. (A.45) to bound Eq. (A.44), with probability at least  $1 - \delta$ , we have

$$\left| \lambda_{\max} \left( \frac{\sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top}{N/2} - \frac{\sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top}{\sum_{i \in [N]} \mathbf{1}_i} \right) \right| \leq \frac{\sqrt{c\ell_\delta N/\ell_1}}{N/2 - \sqrt{c\ell_\delta N/\ell_1}}.$$

Therefore, the first term in Eq. (A.41) is  $o_p(1)$ . It remains to bound  $|\lambda_{\max}(\frac{1}{N} \sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top)|$ . We have

$$\begin{aligned} \left| \lambda_{\max} \left( \frac{1}{N} \sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top \right) \right| &\stackrel{(a)}{=} \left| \max_{x \in \mathbb{R}^N: \|x\|_2=1} x^\top \left( \frac{1}{N} \sum_{i \in [N]} \mathbf{1}_i \tilde{U}_{i,\cdot} \tilde{U}_{i,\cdot}^\top \right) x \right| \\ &= \max_{x \in \mathbb{R}^N: \|x\|_2=1} \frac{1}{N} \sum_{i \in [N]} \mathbf{1}_i |x^\top \tilde{U}_{i,\cdot}|^2 \\ &\stackrel{(b)}{\leq} \max_{x \in \mathbb{R}^N: \|x\|_2=1} \frac{1}{N} \sum_{i \in [N]} \mathbf{1}_i \|x\|_2 \|\tilde{U}_{i,\cdot}\|_2 \stackrel{(c)}{\leq} c, \end{aligned}$$

where (a) follows from the definition of the maximum eigenvalue of a matrix, (b) follows from Cauchy-Schwarz inequality, and (c) follows from Assumption 8(b).

### F.3.2. Assumption 11 holds

First, Assumption 11(a) holds as follows,

$$\max_{j \in [M]} \frac{1}{N} \sum_{j' \in [M]} \left| \sum_{i \in [N]} \mathbb{E}[h_{i,j} h_{i,j'}] \right| \stackrel{(a)}{\leq} \max_{j \in [M]} \frac{1}{N} \sum_{i \in [N]} \sum_{j' \in [M]} \left| \mathbb{E}[h_{i,j} h_{i,j'}] \right| \stackrel{(b)}{\leq} \max_{j \in [M]} \frac{1}{N} \sum_{i \in [N]} c = c,$$

where (a) follows from triangle inequality and (b) follows from Assumption 9(b). Next, from Assumption 9(a) and Assumption 9(c), we have

$$\max_{j \in [M]} |\mathbb{E}[h_{i,j} h_{i',j}]| = \begin{cases} 0 & \text{if } i \neq i' \\ \max_{j \in [M]} |\mathbb{E}[h_{i,j}^2]| \leq c & \text{if } i = i' \end{cases}$$

Then, Assumption 11(b) holds as follows,

$$\max_{i \in [N]} \max_{j \in [M]} \sum_{i' \in [N]} |\mathbb{E}[h_{i,j} h_{i',j}]| \leq c.$$

Next, Assumption 11(c) holds as follows,

$$\frac{1}{NM} \sum_{i, i' \in [N]} \sum_{j, j' \in [M]} |\mathbb{E}[h_{i,j} h_{i',j'}]| \stackrel{(a)}{=} \frac{1}{NM} \sum_{i \in [N]} \sum_{j, j' \in [M]} |\mathbb{E}[h_{i,j} h_{i,j'}]| \stackrel{(b)}{\leq} \frac{1}{NM} \sum_{i \in [N]} \sum_{j \in [M]} c = c,$$

where (a) follows from Assumption 9(c) and (b) follows from Assumption 9(b). Next, let  $\gamma_{i,j,j'} \triangleq h_{i,j} h_{i,j'} - \mathbb{E}[h_{i,j} h_{i,j'}]$  and fix any  $j, j' \in [M]$ . Then, Assumption 11(d) holds as follows,

$$\begin{aligned} \frac{1}{N^2} \mathbb{E} \left[ \left( \sum_{i \in [N]} \gamma_{i,j,j'} \right)^4 \right] &= \frac{1}{N^2} \mathbb{E} \left[ \left( \sum_{i_1 \in [N]} \gamma_{i_1,j,j'} \right) \left( \sum_{i_2 \in [N]} \gamma_{i_2,j,j'} \right) \left( \sum_{i_3 \in [N]} \gamma_{i_3,j,j'} \right) \left( \sum_{i_4 \in [N]} \gamma_{i_4,j,j'} \right) \right] \\ &\stackrel{(a)}{=} \frac{1}{N^2} \sum_{i \in [N]} \mathbb{E} [\gamma_{i,j,j'}^4] + \frac{3}{N^2} \sum_{i \neq i' \in [N]} \mathbb{E} [\gamma_{i,j,j'}^2 \gamma_{i',j,j'}^2] \leq c, \end{aligned}$$

where (a) follows from linearity of expectation and Assumption 9(c) after by noting that  $\mathbb{E}[\gamma_{i,j,j'}] = 0$  for all  $i, j, j' \in [N] \times [M] \times [M]$  and (b) follows because  $\gamma_{i,j,j'}$  has bounded moments due to Assumption 9(a).

## G. Data generating process for the simulations

The inputs of the data generating process (DGP) are: the probability bound  $\lambda$ ; two positive constants  $c^{(0)}$  and  $c^{(1)}$ ; and the standard deviations  $\sigma_{i,j}^{(a)}$  for every  $i \in [N], j \in [M], a \in \{0, 1\}$ . The DGP is:

1. For positive integers  $r_p, r_\theta$  and  $r = \max\{r_p, r_\theta\}$ , generate a proxy for the common unit-level latent factors  $U^{\text{shared}} \in \mathbb{R}^{N \times r}$ , such that, for all  $i \in [N]$  and  $j \in [r]$ ,  $u_{i,j}^{\text{shared}}$  is independently sampled from a  $\text{Uniform}(\sqrt{\lambda}, \sqrt{1-\lambda})$  distribution, with  $\lambda \in (0, 1)$ .
2. Generate proxies for the measurement-level latent factors  $V, V^{(0)}, V^{(1)} \in \mathbb{R}^{M \times r}$ , such that, for all  $i \in [M]$  and  $j \in [r]$ ,  $v_{i,j}, v_{i,j}^{(0)}, v_{i,j}^{(1)}$  are independently sampled from a  $\text{Uniform}(\sqrt{\lambda}, \sqrt{1-\lambda})$  distribution.
3. Generate the treatment assignment probability matrix  $P$

$$P = \frac{1}{r_p} U_{[N] \times [r_p]}^{\text{shared}} V_{[M] \times [r_p]}^\top.$$

4. For  $a \in \{0, 1\}$ , run SVD on  $U^{\text{shared}}V^{(a)\top}$ , i.e.,

$$\text{SVD}(U^{\text{shared}}V^{(a)\top}) = (U^{(a)}, \Sigma^{(a)}, W^{(a)}).$$

Then, generate the mean potential outcome matrices  $\Theta^{(0)}$  and  $\Theta^{(1)}$ :

$$\Theta^{(a)} = \frac{c^{(a)}\text{Sum}(\Sigma^{(a)})}{r_\theta} U_{[N] \times [r_\theta]}^{(a)} W_{[M] \times [r_\theta]}^{(a)\top},$$

where  $\text{Sum}(\Sigma^{(a)})$  denotes the sum of all entries of  $\Sigma^{(a)}$ .

5. Generate the noise matrices  $E^{(0)}$  and  $E^{(1)}$ , such that, for all  $i \in [N], j \in [M], a \in \{0, 1\}$ ,  $\varepsilon_{i,j}^{(a)}$  is independently sampled from a  $\mathcal{N}(0, (\sigma_{i,j}^{(a)})^2)$  distribution. Then, determine  $y_{i,j}^{(a)}$  from Eq. (2).
6. Generate the noise matrix  $W$ , such that, for all  $i \in [N], j \in [M]$ ,  $\eta_{i,j}$  is independently sampled as per Eq. (4). Then, determine  $a_{i,j}$  and  $y_{i,j}$  from Eq. (3) and Eq. (1), respectively.

In our simulations, we set  $\lambda = 0.05$ ,  $c^{(0)} = 1$  and  $c^{(1)} = 2$ . In practice, instead of choosing the values of  $\sigma_{i,j}^{(a)}$  as ex-ante inputs, we make them equal to the standard deviation of all the entries in  $\Theta^{(a)}$  for every  $i$  and  $j$ , separately for  $a \in \{0, 1\}$ .

## References

- Abadie, A., Diamond, A., and Hainmueller, J. (2010). Synthetic control methods for comparative case studies: Estimating the effect of california’s tobacco control program. *Journal of the American Statistical Association*, 105(490):493–505.
- Abadie, A. and Imbens, G. W. (2006). Large sample properties of matching estimators for average treatment effects. *Econometrica*, 74(1):235–267.
- Agarwal, A., Dahleh, M., Shah, D., and Shen, D. (2023a). Causal matrix completion. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 3821–3826. PMLR.
- Agarwal, A., Shah, D., and Shen, D. (2023b). Synthetic interventions.
- Agarwal, A., Shah, D., Shen, D., and Song, D. (2021). On robustness of principal component regression. *Journal of the American Statistical Association*, pages 1–34.
- Angrist, J. D. (1998). Estimating the labor market impact of voluntary military service using social security data on military applicants. *Econometrica*, 66(2):249–288.
- Arkhangelsky, D., Athey, S., Hirshberg, D. A., Imbens, G. W., and Wager, S. (2021). Synthetic difference-in-differences. *American Economic Review*, 111(12):4088–4118.

- Athey, S., Bayati, M., Doudchenko, N., Imbens, G., and Khosravi, K. (2021). Matrix completion methods for causal panel data models. *Journal of the American Statistical Association*, 116(536):1716–1730.
- Bai, J. (2003). Inferential theory for factor models of large dimensions. *Econometrica*, 71(1):135–171.
- Bai, J. (2009). Panel data models with interactive fixed effects. *Econometrica*, 77(4):1229–1279.
- Bai, J. and Ng, S. (2002). Determining the number of factors in approximate factor models. *Econometrica*, 70(1):191–221.
- Bai, J. and Ng, S. (2021). Matrix completion, counterfactuals, and factor analysis of missing data. *Journal of the American Statistical Association*, 116(536):1746–1763.
- Bang, H. and Robins, J. M. (2005). Doubly robust estimation in missing data and causal inference models. *Biometrics*, 61(4):962–972.
- Bhatia, R. (2007). *Perturbation bounds for matrix eigenvalues*. SIAM.
- Bhattacharya, S. and Chatterjee, S. (2022). Matrix completion with data-dependent missingness probabilities. *IEEE Transactions on Information Theory*, 68(10):6762–6773.
- Billingsley, P. (2017). *Probability and measure*. John Wiley & Sons.
- Chatterjee, S. (2015). Matrix estimation by universal singular value thresholding. *The Annals of Statistics*, 43(1):177 – 214.
- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., and Robins, J. (2018). Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68.
- Cochran, W. G. (1968). The effectiveness of adjustment by subclassification in removing bias in observational studies. *Biometrics*, pages 295–313.
- Dwivedi, R., Tian, K., Tomkins, S., Klasnja, P., Murphy, S., and Shah, D. (2022a). Counterfactual inference for sequential experiments. *arXiv preprint arXiv:2202.06891*.
- Dwivedi, R., Tian, K., Tomkins, S., Klasnja, P., Murphy, S., and Shah, D. (2022b). Doubly robust nearest neighbors in factor models. *arXiv preprint arXiv:2211.14297*.
- Imbens, G. W. and Rubin, D. B. (2015). *Causal Inference for Statistics, Social, and Biomedical Sciences: An Introduction*. Cambridge University Press.
- Li, Y., Shah, D., Song, D., and Yu, C. L. (2019). Nearest neighbors for matrix estimation interpreted as blind regression for latent variable model. *IEEE Transactions on Information Theory*, 66(3):1760–1784.

- Ma, W. and Chen, G. H. (2019). Missing not at random in matrix completion: The effectiveness of estimating missingness probabilities under a low nuclear norm assumption. *Advances in neural information processing systems*, 32.
- Nguyen, L. T., Kim, J., and Shim, B. (2019). Low-rank matrix completion: A contemporary survey. *IEEE Access*, 7:94215–94237.
- Rosenbaum, P. R. and Rubin, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55.
- Tropp, J. A. (2012). User-friendly tail bounds for sums of random matrices. *Foundations of computational mathematics*, 12:389–434.
- Vershynin, R. (2018). *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press.