

## INFERENCE FOR LOW-RANK MODELS

BY VICTOR CHERNOZHUKOV<sup>1,a</sup>, CHRISTIAN HANSEN<sup>2,b</sup>, YUAN LIAO<sup>3,c</sup> AND  
YINCHU ZHU<sup>4,d</sup>

<sup>1</sup>*Department of Economics, Massachusetts Institute of Technology, [vchern@mit.edu](mailto:vchern@mit.edu)*

<sup>2</sup>*Booth School of Business, University of Chicago, [Christian.Hansen@chicagobooth.edu](mailto:Christian.Hansen@chicagobooth.edu)*

<sup>3</sup>*Department of Economics, Rutgers University, [yuan.liao@rutgers.edu](mailto:yuan.liao@rutgers.edu)*

<sup>4</sup>*Department of Economics, Brandeis University, [yinchuzhu@brandeis.edu](mailto:yinchuzhu@brandeis.edu)*

This paper studies inference in linear models with a high-dimensional parameter matrix that can be well approximated by a “spiked low-rank matrix.” A spiked low-rank matrix has rank that grows slowly compared to its dimensions and nonzero singular values that diverge to infinity. We show that this framework covers a broad class of models of latent variables, which can accommodate matrix completion problems, factor models, varying coefficient models and heterogeneous treatment effects. For inference, we apply a procedure that relies on an initial nuclear-norm penalized estimation step followed by two ordinary least squares regressions. We consider the framework of estimating incoherent eigenvectors and use a rotation argument to argue that the eigenspace estimation is asymptotically unbiased. Using this framework, we show that our procedure provides asymptotically normal inference and achieves the semiparametric efficiency bound. We illustrate our framework by providing low-level conditions for its application in a treatment effects context where treatment assignment might be strongly dependent.

**1. Introduction.** We study inference for linear low-rank models:

$$Y = X \circ \Theta + \mathcal{E},$$

where  $(Y, X, \Theta, \mathcal{E})$  are  $n \times p$  matrices with both  $n, p \rightarrow \infty$  and  $\circ$  denotes the matrix element-wise product. We observe data  $(X, Y)$ , and  $\mathcal{E}$  represents unobserved statistical noise. The model parameter is the matrix coefficient  $\Theta$ . We assume  $\Theta$  follows an *approximate spiked low-rank model*:  $\Theta$  can be well approximated by a low-rank matrix whose rank  $J$  is either fixed or grows slowly compared to  $n, p$  and whose largest  $J$  singular values diverge with  $(n, p)$ . Our main goal is performing statistical inference on both sparse and dense linear combinations of elements of  $\Theta$ .

Under the approximate spiked low-rank model structure, nuclear-norm regularization provides a natural benchmark approach to estimating  $\Theta$ . There is a substantial literature that studies rates of convergence of nuclear-norm penalized estimators; see, for example, [Koltchinskii, Lounici and Tsybakov \(2011\)](#) and [Negahban and Wainwright \(2011\)](#) for prominent examples. Providing results in low-rank models that allow for ready construction of inferential objects such as confidence intervals has been a topic in the more recent literature. For example, [Xia and Yuan \(2021\)](#) and [Chen et al. \(2019\)](#) study inference in settings where the matrix parameter of interest has an exact low-rank structure with fixed rank and elements of  $X$  are i.i.d. copies from an unknown distribution.

We contribute to this literature by establishing asymptotic normality for low-rank estimators. Our method starts with an initial estimator of  $\Theta$  obtained using nuclear-norm regularization from which we extract the right singular vectors. We then treat the extracted singular

vectors as data and obtain estimates of the left singular vectors and updated estimates of the right singular vectors by applying additional least squares steps. The final estimator is then the product of the estimated left and right singular vectors. We make use of a rotation argument to show that, in terms of estimating the space of the singular vectors, the regularization bias of the first step nuclear-norm penalized estimation aligns with the space spanned by the true singular vectors. Thus, the inference using additional least squares steps is not affected by the regularization bias.

We prove that our estimator for linear functionals of the low-rank matrix is asymptotically normal. We also establish the semiparametric efficiency bound and show that our estimator attains the efficiency bound. The notion of semiparametric efficiency in the presence of high-dimensional nuisance parameters is adopted from [Janková and van de Geer \(2018\)](#). Our result is novel relative to [Janková and van de Geer \(2018\)](#) because they deal with sparse models while we look at linear combinations of a low-rank matrix.

Our conditions allow possible strong dependence within  $X$ , which is useful in many contexts. For example, in the matrix completion context, we can accommodate persistence in observed and missing entries rather than relying on independent missingness. In our treatment effects example, allowing strong dependence allows us to consider scenarios where units are first in the control state for a period of time and then enter the treated state and remain there until the end of the sample period.

We rely on two key technical conditions in establishing asymptotic normality of our proposed estimator. We first assume  $\Theta$  has *spiked singular values* (SSV), which requires that the nonzero singular values are large. This condition ensures that the rank of  $\Theta$  can be consistently estimated and that the singular vectors are estimated sufficiently well for use in Stage 2 of the procedure. In the inference context, to make *entrywise* inference for the low-rank matrix, the SSV condition on the singular values seem necessary.

The second condition relates to *incoherent singular vectors* as defined in, for example, [Candès and Plan \(2010\)](#), [Candès and Recht \(2009\)](#), [Keshavan, Montanari and Oh \(2010\)](#) and [Chen et al. \(2020\)](#). This condition requires that the signals on the singular vectors should be approximately evenly distributed across their entries. Under the incoherence condition, we use a “rotation” argument to show that our approach provides asymptotically unbiased estimates of the eigenvector space. We note that incoherence does rule out the setting of “sparse PCA,” which needs a separate treatment and often requires explicit debiasing steps as in, for example, [Janková and van de Geer \(2021\)](#).

We note that the SSV and incoherence conditions are strong and are often absent in the literature when probability bounds are derived. However, probability bounds for the Frobenius risk in general cannot imply the asymptotic distribution of estimators. In particular, one of the objects of interest in this paper is to make inference on sparse linear combinations of rows (or columns) of  $\Theta$ , including elementwise inference. Recent developments for perturbation bounds of entrywise eigenanalysis require SSV to make *entrywise* inference for a low-rank matrix; see, for example, [Abbe et al. \(2020\)](#).

To further illustrate that both conditions seem necessary for good performance, we provide new minimax theory on convergence rates without them. These results verify that it is impossible to guarantee entrywise consistency without SSV or incoherence and show that, for dense linear combinations, the optimal rates one can achieve without these conditions are potentially much worse than those available under them. Finally, as these minimax rates do show that inference for dense linear combinations may proceed without SSV or incoherence, we provide an alternative inference for dense functionals when these conditions are relaxed.

The low-rank model being considered has wide applicability. We show that the function class of reproducing kernel Hilbert space (RKHS) can be approximated using a spiked low-rank model, and we verify the SSV and incoherent conditions in RKHS.

*The literature.* Low-rank regression has been extensively studied in the literature. Much of this work focuses on deriving sharp deviation bounds for low-rank estimators; see, for instance, Recht (2011), Gross (2011), Rohde and Tsybakov (2011), Koltchinskii, Lounici and Tsybakov (2011), Dray and Josse (2015), Zhu, Wang and Samworth (2022), Candès and Plan (2010), Hastie et al. (2015), Keshavan, Montanari and Oh (2010) and Sun and Zhang (2012). As with Xia and Yuan (2021) and Chen et al. (2019), our paper complements the literature by providing asymptotic distributional results. A key difference between our work and Xia and Yuan (2021) and Chen et al. (2019) is that our approach does not rely on explicit debiasing steps to achieve asymptotic normality. Rather, we rely on the fact that  $\Theta$  is a “product parameter” obtained by multiplication of left and right singular vectors to verify that our procedure produces sufficiently regular estimators for asymptotic normality to hold without explicit debiasing.

Our paper is related to Chernozhukov et al. (2018), which considers inference in linear panel data models with multivariate coefficient matrices that admit factor structures. There are several important differences between the two papers. Because Chernozhukov et al. (2018) consider estimation of multiple matrix parameters, they employ a complicated orthogonalization step to deal with the fact that regularization bias in any of the matrix parameters spills over and impacts estimation of all other matrix parameters. Chernozhukov et al. (2018) also rely on strong conditions on regressors while our conditions allow the regressor in our model to be strongly persistent. This generalization allows us to handle matrix completion problems with “systematic missingness”. We also provide several new optimality results to establish semiparametric efficiency and minimaxity. Finally, we explicitly allow for the low-rank structure to be an approximation by accounting for approximation errors and allowing the rank of the approximating low-rank structure to increase with the sample size. Admitting these characteristics broadens the applicability of the method. For example, nonparametric models under RKHS cannot be formulated as exact low-rank models with fixed rank but can be approximated by low-rank models with slowly growing rank.

Throughout the paper, we denote the maximum and minimum singular values of a matrix  $A$  as  $\psi_{\max}(A)$  and  $\psi_{\min}(A)$ . We use  $\psi_j(A)$  to denote the  $j$ th largest singular value of  $A$ . We use  $\|A\|_F$ ,  $\|A\|$  and  $\|A\|_{(n)} = \sum_{k=1}^{\min\{n,p\}} \psi_k(A)$  to respectively denote the matrix Frobenius norm, operator norm and nuclear norm. We let  $\|A\|_{\max} = \max_{i,j} |(A)_{ij}|$  be the elementwise norm. Let  $\text{vec}(A)$  denote the vector that stacks the columns of  $A$ . For two stochastic sequences, we write  $a_n \asymp b_n$  if  $a_n = O_P(b_n)$  and  $b_n = O_P(a_n)$ , which means  $a_n/b_n = O_P(1)$ . Finally,  $a \vee b$  means  $\max(a, b)$ .

**2. Estimation procedure.**

2.1. *Spiked low-rank matrices.* Consider the following model:

$$(1) \quad y_{ij} = x_{ij}\theta_{ij} + \varepsilon_{ij}, \quad i \leq n, j \leq p,$$

where we observe data  $(y_{ij}, x_{ij})$  and  $\varepsilon_{ij}$  is the noise term. Let  $(Y, X, \Theta, \mathcal{E})$  denote the  $n \times p$  matrices of  $(y_{ij}, x_{ij}, \theta_{ij}, \varepsilon_{ij})$ . Then the matrix form of (1) is

$$Y = X \circ \Theta + \mathcal{E},$$

where  $\circ$  denotes the matrix elementwise product. The goal is to make inference about linear combinations of elements of  $\Theta$ . Throughout the paper, we impose that  $\Theta$  and its associated singular values/vectors are random. Suppose  $\Theta$  can be decomposed as

$$(2) \quad \Theta = \Theta_0 + R,$$

where  $\Theta_0$  and  $R$  are  $n \times p$  matrices satisfying the following conditions:

(i)  $\Theta_0$  is a rank  $J$  matrix where  $J$  is either bounded or grows slowly compared to  $(n, p)$ . In addition, the nonzero singular values of  $\Theta_0$  are “spiked”:

$$\psi_1(\Theta_0) \geq \dots \geq \psi_J(\Theta_0) \geq \psi_{np}, \quad \psi_j(\Theta_0) = 0 \quad \forall j > J$$

for some sequence  $\psi_{np} \rightarrow \infty$ .

(ii)  $R$  is the low-rank approximation error whose entries  $r_{ij}$  satisfy

$$\max_{i,j} |r_{ij}| \leq O_P(r_{np})$$

for some sequence  $r_{np} \rightarrow 0$ .

(iii) Let  $U_0 = [u_1, \dots, u_n]'$  and  $V_0 = [v_1, \dots, v_p]'$ , respectively, denote the  $n \times J$  and  $p \times J$  matrices that collect the left singular vectors and right singular vectors of  $\Theta_0$  corresponding to the nonzero singular values. We assume incoherent singular-vectors:

$$\max_{j \leq p} \|v_j\| = O_P(\sqrt{Jp^{-1}}), \quad \max_{i \leq n} \|u_i\| = O_P(\sqrt{Jn^{-1}}).$$

Given the approximate low-rank structure of  $\Theta$ , a natural estimation strategy is nuclear-norm penalized optimization:

$$(3) \quad \tilde{\Theta} = \arg \min_{\Theta \in \mathcal{A}} \|Y - X \circ \Theta\|_F^2 + \nu \|\Theta\|_{(n)},$$

where  $\mathcal{A} = \{\Theta : \|\Theta\|_{\max} \leq M\}$  and  $\nu$  is a tuning parameter. Imposing the max-norm constraint with a large constant  $M > 0$  helps stabilize the solution; see, for example, [Klopp \(2014\)](#). Nuclear-norm penalized regression is natural as the solution is easy to compute. Another option would be to explicitly penalize the matrix rank; however, obtaining the solution to the rank penalized problem is in general difficult unless all elements of  $X$  are equal to one. Statistical properties of (3), focusing on the minimax rate for  $\|\tilde{\Theta} - \Theta\|_F$ , have been well studied in the literature; see, for example, [Koltchinskii, Lounici and Tsybakov \(2011\)](#) and [Negahban and Wainwright \(2011\)](#). It is also well known that the singular values of  $\tilde{\Theta}$  suffer from shrinkage biases, so  $\tilde{\Theta}$  is not suitable for inference.

We assume that  $J$ , the rank of the low-rank component  $\Theta_0$ , is known for simplicity. For instance, in the treatment effect study where the parameter matrix is approximated by a low-rank structure via a sieve representation, the rank equals the sieve dimension, which could be prespecified. In cases where rank is unknown, it can be consistently estimated. For example, one can apply the singular value thresholding method where the cut-off value for “large singular values” can be chosen to dominate the noise level; see, for example, [Onatski \(2010\)](#) and [Fan, Guo and Zheng \(2022\)](#).

2.2. *The proposed estimation procedure.* Let the singular value decomposition of  $\Theta_0$  be

$$\Theta_0 = U_0 D_0 V_0' := \Gamma_0 V_0', \quad \Gamma_0 := U_0 D_0.$$

Here,  $D_0$  is the  $J \times J$  diagonal matrix containing the nonzero singular values of  $\Theta_0$ , and  $U_0$  and  $V_0$  are respectively the  $n \times J$  left singular vector matrix of  $\Theta_0$  and  $p \times J$  right singular vector matrix of  $\Theta_0$  corresponding to the nonzero singular values. Let  $\gamma'_{0,i}$  for  $i = 1, \dots, n$  denote the rows of  $\Gamma_0$ , and let  $v'_{0,j}$  for  $j = 1, \dots, p$  denote the rows of  $V_0$ .

ALGORITHM 2.1. Fix  $i \leq n$ . Estimate  $\theta_{ij}$  ( $j = 1, \dots, p$ ) as follows:

Step 1 *Sample splitting.* Randomly split the sample into  $\{1, \dots, n\} \setminus \{i\} = \mathcal{I} \cup \mathcal{I}^c$  disjointly, so that  $|\mathcal{I}|_0 = \lfloor (n - 1)/2 \rfloor$ . Let

$$\mathcal{G}_{\mathcal{I}} := (Y_{\mathcal{I}}, X_{\mathcal{I}}, \Theta_{\mathcal{I}}),$$

respectively, denote the  $|\mathcal{I}|_0 \times p$  submatrices of  $\mathcal{G} := (Y, X, \Theta)$  for observations  $i \in \mathcal{I}$ . Estimate the low-rank matrix  $\Theta_{\mathcal{I}}$  as

$$(4) \quad \tilde{\Theta}_{\mathcal{I}} = \arg \min_{\|\Theta_{\mathcal{I}}\|_{\max} < M} \|Y_{\mathcal{I}} - X_{\mathcal{I}} \circ \Theta_{\mathcal{I}}\|_F^2 + \nu \|\Theta_{\mathcal{I}}\|_{(n)}.$$

We provide a specific feasible choice for  $\nu$  when discussing the simulation example in Section 6. Let  $\tilde{V}_{\mathcal{I}} = (\tilde{v}_1, \dots, \tilde{v}_p)'$  be the  $p \times J$  matrix whose columns are the first  $J$  eigenvectors of  $\tilde{\Theta}'_{\mathcal{I}} \tilde{\Theta}_{\mathcal{I}}$ .

Step 2 *Unbiased estimate of  $\Gamma_0, V_0$ .* Using data  $\mathcal{I}^c$ , obtain

$$\hat{\gamma}_{k,\mathcal{I}} = \arg \min_{\gamma} \sum_{j=1}^p [y_{kj} - x_{kj} \cdot \gamma' \tilde{v}_j]^2, \quad k \in \mathcal{I}^c \cup \{i\}.$$

Update estimates of  $V_0$  as  $\hat{V}_{\mathcal{I}} = (\hat{v}_{1,\mathcal{I}}, \dots, \hat{v}_{p,\mathcal{I}})'$ , where

$$\hat{v}_{j,\mathcal{I}} = \arg \min_v \sum_{k \in \mathcal{I}^c \cup \{i\}} [y_{kj} - x_{kj} \cdot \hat{\gamma}'_{k,\mathcal{I}} v]^2, \quad j = 1, \dots, p.$$

Step 3 *Exchange  $\mathcal{I}$  and  $\mathcal{I}^c$ .* Repeat steps 1–2 with  $\mathcal{I}$  and  $\mathcal{I}^c$  exchanged to obtain  $\hat{\gamma}_{k,\mathcal{I}^c}$  for  $k \in \mathcal{I} \cup \{i\}$  and  $\hat{V}_{\mathcal{I}^c}$ . Define the estimator for  $\theta_{ij}$  as

$$\hat{\theta}_{ij} = \frac{1}{2} [\hat{\gamma}'_{i,\mathcal{I}} \hat{v}_{j,\mathcal{I}} + \hat{\gamma}'_{i,\mathcal{I}^c} \hat{v}_{j,\mathcal{I}^c}].$$

We only iterate least squares *once* in step 2. The least squares steps following the use of nuclear-norm penalized estimation are analogous to approaches in the sparse regression setting that rely on refitting the least squares using selected regressors in a first step, such as post-lasso, for example, [Belloni and Chernozhukov \(2013\)](#). The motivation is similar in wanting to alleviate shrinkage biases induced in the initial penalized estimation step. In addition, we split the sample  $\{1, \dots, n\} \setminus \{i\} = \mathcal{I} \cup \mathcal{I}^c$  so that  $i$  is excluded from both subsamples. Splitting in this way ensures that the  $\varepsilon_{ij}$  for the  $i$  of interest are independent of observations in both subsamples assuming independence across  $i$ .

Stage 2, which involves two least squares estimation steps, is the essential stage to alleviating shrinkage bias. It starts with treating  $\tilde{V}$  from the penalized regression as observed data. A key ingredient of the analysis is to establish that this step produces an approximately unbiased estimator  $\hat{\Gamma}$ , which then allows construction of a well-behaved estimator  $\hat{\Theta}$  in the final step. Given its importance, we provide the intuition for this step in Section 3.

REMARK 2.1. The proposed procedure is similar to the “alternating minimization” (AltMin) method in the literature, for example, [Hastie et al. \(2015\)](#) and [Jain, Netrapalli and Sanghavi \(2013\)](#). There are two key differences. The first is that the AltMin procedure would iterate until convergence. In contrast, we only iterate once and good asymptotic statistical properties are guaranteed. The second difference is that penalization is often carried throughout iterations in the AltMin procedure. Thus, AltMin-type estimators have asymptotic shrinkage biases, which complicates establishing asymptotic normality. By employing unpenalized least squares in Stage 3, our procedure ensures the final estimator does not have large shrinkage bias asymptotically.

**3. Discussion.** We make use of a “rotation” argument and the structure of the low-rank matrix parameter to prove that eigenspace estimation is approximately unbiased if singular vectors are incoherent. Before turning to the matrix parameter setting, we introduce the main idea in the context of estimating a scalar parameter that is itself a product of two parameters.

3.1. *Inference about product parameters.* Consider the problem of estimating a scalar parameter  $\theta$  that can be written as the product of another two scalar parameters:

$$\theta = \gamma\beta, \quad \gamma, \beta \in \mathbb{R}.$$

Suppose some initial estimate  $\tilde{\beta}$  can be obtained for  $\beta$ , which is consistent but may have first-order bias. In addition, suppose that  $\theta$  can be identified as the unique minimizer of a population loss function:

$$\theta = \arg \min_{a \in \mathcal{A}} Q(a),$$

where  $Q(\cdot)$  is the loss function and  $\mathcal{A}$  is the parameter space. Let  $Q_n(\cdot)$  denote the sample version of  $Q$  and suppose both  $Q_n$  and  $Q$  are twice continuously differentiable. Let  $\dot{Q}_n(a) = \frac{d}{da} Q_n(a)$ ,  $\dot{Q}(a) = \frac{d}{da} Q(a)$  and  $\ddot{Q}(a) = \frac{d^2}{da^2} Q(a)$ . Let  $(\gamma, \beta, \theta)$  represent the true values of the parameters.

We consider an iterative procedure to estimate  $\theta$  that mimics the approach we propose in the matrix parameter setting:

- (i) Obtain  $\hat{\gamma} = \arg \min_{\gamma} Q_n(\gamma\tilde{\beta})$
- (ii) Obtain  $\hat{\beta} = \arg \min_{\beta} Q_n(\hat{\gamma}\beta)$
- (iii) Set  $\hat{\theta} = \hat{\gamma}\hat{\beta}$ .

In step (i), standard analysis based on Taylor expansion leads to

$$(5) \quad \hat{\gamma} - \gamma = G^{-1}\beta\dot{Q}_n(\theta) + G^{-1}\partial_{\gamma,\beta}^2 Q(\gamma\beta)(\tilde{\beta} - \beta) + o(|\hat{\gamma} - \gamma|),^1$$

where  $G = -\partial_{\gamma,\gamma}^2 Q_n(\gamma\beta)$ . The first term in the expansion is the score, which leads to asymptotic normality in usual cases. The second term reflects the effect of the initial estimate  $\tilde{\beta}$ .

In general, the second term will lead to poor performance of  $\hat{\gamma}$  if the initial estimator  $\tilde{\beta}$  is ill-behaved. One approach, dating back to at least [Neyman \(1959\)](#), is to rely on appropriately “orthogonalized” scores. This property would correspond to basing estimation on an objective function that satisfied  $\partial_{\gamma,\beta}^2 Q(\gamma\beta) = 0$  at the population level in the present case; see, for example, [Chernozhukov, Hansen and Spindler \(2015\)](#) for a review of such approaches.

The fact that the “product parameter”  $\theta$ , rather than  $\gamma$  itself, is the object of interest allows a new argument in this paper. The key is that the loss function depends on  $\theta$  only through the product of  $(\gamma, \beta)$ . It is straightforward to verify that

$$\partial_{\gamma,\beta}^2 Q(\gamma\beta) = \gamma \underbrace{\ddot{Q}(\theta)}_{\text{score}=0} \beta + \dot{Q}(\theta) = \gamma \ddot{Q}(\theta) \beta.$$

Substituting this expression for  $\partial_{\gamma,\beta}^2 Q(\gamma\beta)$  into (5) then produces

$$\hat{\gamma} - \gamma = G^{-1}\beta\dot{Q}_n(\theta) + G^{-1}\gamma\ddot{Q}(\theta)\beta(\tilde{\beta} - \beta) + o(|\hat{\gamma} - \gamma|).$$

An important observation is that the second term  $G^{-1}\gamma\ddot{Q}(\theta)\beta(\tilde{\beta} - \beta)$  is proportional to  $\gamma$ . We can move it to the left-hand side of the expansion for  $\hat{\gamma}$  to obtain

$$\hat{\gamma} - H\gamma = G^{-1}\beta\dot{Q}_n(\theta) + o(|\hat{\gamma} - \gamma|)$$

for  $H := 1 + G^{-1}\ddot{Q}(\theta)\beta(\tilde{\beta} - \beta)$ . Hence,  $\hat{\gamma}$  estimates a “rotated” version of  $\gamma$  with no first-order bias. As such, in the sense of estimating the “space” of  $\gamma$ , the effect  $\tilde{\beta} - \beta$  is negligible

<sup>1</sup>We have  $0 = \partial_{\gamma} Q_n(\hat{\gamma}\tilde{\beta}) = \partial_{\gamma} Q_n(\gamma\beta) + \partial_{\gamma,\beta}^2 Q_n(\gamma\beta)(\tilde{\beta} - \beta) + \partial_{\gamma,\gamma}^2 Q_n(\gamma\beta)(\hat{\gamma} - \gamma) + O(|\hat{\gamma} - \gamma|^2 + |\tilde{\beta} - \beta|^2)$ , and  $\partial_{\gamma} Q_n(\gamma\beta) = \beta\dot{Q}_n(\theta)$ . Inverting  $\partial_{\gamma,\gamma}^2 Q_n(\gamma\beta)$  leads to (5).

as it is “absorbed” by the rotation matrix. In addition,  $H$  is asymptotically invertible since  $H \xrightarrow{P} 1$ .

Moving on to step (ii), it is clear that  $\widehat{\beta}$  estimated in this step will be an approximately unbiased estimator for  $H^{-1}\beta$ . The rotation matrices will then cancel in estimating the parameter of interest:

$$\widehat{\theta} := \widehat{\gamma}\widehat{\beta} = \gamma H H^{-1}\beta + o_P(1) = \theta + o_P(1).$$

After appropriate scaling, the leading term hidden in the  $o_P(1)$  in the final expression will also be asymptotically normal. It is this cancellation of rotation matrices that underlies our “rotation-unbiasedness.” Furthermore, in models where  $\sqrt{n}$ -consistency is attainable,  $\sqrt{n}(\widehat{\theta} - \theta)$  is asymptotically normal as long as the initial estimator satisfies  $|\widetilde{\beta} - \beta| = o_P(n^{-1/4})$ .

The intuition of “rotation-unbiasedness” as described above has also been observed previously in the literature. [Keshavan, Montanari and Oh \(2010\)](#) studied local geometric properties in Grassmann manifold and related optimization algorithms. [Sun and Luo \(2016\)](#) examined the local geometry of the loss  $f(\Gamma, V) = \|Y - \Gamma V'\|_F^2$  in the matrix completion context. Our observation aligns with theirs, but we use this observation in the context of estimation bias. In our setting, the geometry of product-parameter  $\gamma\beta$  ensures that the effect of first-step estimation error  $\widetilde{\beta} - \beta$  is aligned with the space of the true  $\gamma$ . This alignment results in our ability to establish asymptotic normality of our final estimator without relying on any additional debiasing schemes beyond the use of a single set of least squares steps in step 2 of our algorithm.

**3.2. Eigenspace estimation.** In the low-rank inference context, recall that  $\Theta_0 = \Gamma_0 V_0'$ , which is the product of two parameters. Related to the simple example in the previous section, we think about  $V_0$  as  $\beta$  and use the singular vectors  $\widetilde{V}$  extracted from the nuclear-norm regularized estimator as its initial estimate.

Write  $\widehat{\Gamma} = (\widehat{\gamma}_1, \dots, \widehat{\gamma}_n)'$  and  $\widetilde{V} = (\widetilde{v}_1, \dots, \widetilde{v}_p)'$ . Then for each  $i \leq n$ ,

$$\widehat{\gamma}_i = \arg \min_{\gamma} Q_i(\gamma, \widetilde{V}), \quad Q_i(\gamma, \widetilde{V}) := \sum_{j=1}^p [y_{ij} - x_{ij} \cdot \gamma' \widetilde{v}_j]^2.$$

Then for some  $J \times J$  matrix  $G^{-1}$ , Taylor expansion leads to

$$\widehat{\gamma}_i - \gamma_i = G^{-1} \partial_{\gamma} Q_i(\gamma_i, V_0) + \frac{\partial^2 Q_i(\gamma_i, V_0)}{\partial \gamma \partial \text{vec}(V)} \text{vec}(\widetilde{V} - V_0) + \text{higher-order terms}.$$

The leading term  $G^{-1} \partial_{\gamma} Q_i(\gamma_i, V)$  is asymptotically normal if  $V_0$  is incoherent. The second term satisfies

$$\frac{\partial^2 Q_i(\gamma_i, V)}{\partial \gamma \partial \text{vec}(V)} \text{vec}(\widetilde{V} - V_0) = H_1 \gamma_i + \Delta_i$$

for some rotation matrix  $H_1$  and higher-order term  $\Delta_i$ .

The term  $H_1 \gamma_i$  is a rotated version of  $\gamma_i$ . Defining  $H := I + H_1$  and moving  $H_1 \gamma_i$  to the left-hand side then yields the matrix form expansion:

$$\widehat{\Gamma} - \Gamma_0 H = \partial_{\Gamma} Q_p(\Gamma, V_0) G^{-1} + \text{higher-order terms},$$

where  $\partial_{\Gamma} Q_p(\Gamma, V_0)$  is an  $n \times J$  matrix whose  $i^{\text{th}}$  row is the transpose of  $\partial_{\gamma} Q_i(\gamma_i, V_0)$ . Following the logic outlined in Section 3.1, we have that the follow-up estimator  $\widehat{V}$  will recover an appropriately rotated version of  $V$  to cancel with  $H$ . Consequently,

$$\widehat{\Theta} = \widehat{\Gamma} \widehat{V}' \approx \Gamma_0 H H^{-1} V_0' = \Gamma_0 V_0' = \Theta_0.$$

It will then follow that  $\widehat{\Theta}$  is approximately unbiased with sampling distribution that can be approximated by a centered Gaussian distribution. As in the simpler scalar case, the key feature we take advantage of is that we only need the estimated  $V_0$  to have the same span as the actual  $V_0$  if our goal is inference about  $\Theta$  or the space spanned by the singular vectors.

3.3. *Sample splitting.* Our argument for demonstrating that the higher-order term,  $\Delta_i$  is asymptotically negligible relies on sample splitting. The structure of  $\Delta_i$  is

$$\Delta_i = B \sum_{j=1}^p (\tilde{v}_j - v_j) \varepsilon_{ij} x_{ij}$$

for some matrix  $B$ .

For a fixed  $i$ , let  $\mathcal{I} \subset \{1, \dots, n\} \setminus i$  be a subset of unit indexes that does not include  $i$ ; and let

$$\mathcal{D}_{\mathcal{I}} = \{(y_{kj}, x_{kj}) : k \in \mathcal{I}, j \leq p\}.$$

Our approach uses only data  $\mathcal{D}_{\mathcal{I}}$ , rather than making use of the full data set, for the initial nuclear-norm penalized regression from which we extract singular vectors for the subsequent OLS rotation-debiasing step. Maintaining independence across  $i$ , estimation errors in the initial estimator of the singular vectors are then independent of variables indexed by  $i$  because  $i \notin \mathcal{I}$ . Assuming  $\varepsilon_{ij}$  is independent across subjects  $i = 1, \dots, n$ , we then have that  $\varepsilon_{ij} x_{ij}$  is independent of estimation error in the singular vectors,  $\tilde{v}_j - v_j$ . We can then easily argue that  $\Delta_i$  has no impact on the asymptotic distribution of the final estimator.

**4. Asymptotic results.** We now present our main results. In Section 4.1, we lay out key conditions and state our result on asymptotic normality. We then provide a brief discussion of semiparametric efficiency in Section 4.2 and then highlight the role of the key SSV and incoherence conditions in Section 4.3 where we present novel minimax results. Finally, we present an alternative estimation scheme for dense linear combinations in Section 4.4.

4.1. *Asymptotic normality.* The goal is to establish inferential theory for the linear functional  $\theta'_i g$ . Here,  $\theta'_i$  denotes the  $i$ th row of  $\Theta$ , and  $g = (g_1, \dots, g_p)' \in \mathbb{R}^p$  is a vector of weights of interest with nonzero weights collected in

$$\mathcal{G} = \{j \leq p : g_j \neq 0\}.$$

Inference on a linear combination of a column of  $\Theta$  can be carried out similarly by switching the roles of  $i$  and  $j$ . Two examples of  $g$  are of particular interest.

*Sparse weights:*  $g$  is a sparse vector with a bounded number of nonzero elements:

$$(6) \quad |\mathcal{G}| = O(1).$$

$\theta'_i g$  thus corresponds to a linear combination of a small number of elements and may be used when we are particularly interested in just a few components of  $\theta_i$ . The sparse  $g$  scenario includes  $g = e_j$  where  $e_j = (0, \dots, 0, 1, 0, \dots, 0)$  is the  $j$ th standard vector for a particular  $j$  in which case  $\theta'_i g = \theta_{ij}$ .

*Dense weights:*  $g$  is a dense vector, in the sense that  $|\mathcal{G}| = O(p)$ , but

$$(7) \quad \max_{j \leq p} |g_j| < Cp^{-1} \quad \text{for some } C > 0.$$

In this case,  $\theta'_i g$  typically represents a weighted average of all components of  $\theta_i$  and includes  $g = (\frac{1}{p}, \dots, \frac{1}{p})'$  as a special case.

The following assumption formally quantifies the requirement of  $g$ . Consider the matrix of standardized right singular vectors:

$$\bar{V}' = \sqrt{p} V_0'.$$



ASSUMPTION 4.1. For some constants  $c, C > 0$ ,

$$c < \|\bar{V}'g\| \leq C, \quad \|g\| < C.$$

In addition,  $g$  satisfies either (6) or (7).

The next assumption restricts the noise data generating process (DGP).

ASSUMPTION 4.2 (DGP for  $\varepsilon_{ij}$ ). (i)  $\varepsilon_{ij}$  is conditionally independent across  $i \leq n$  and  $j \leq p$ , given  $(\Theta, X)$ . Also,  $E(\varepsilon_{ij}|\Theta, X) = 0$  and  $\max_{ij} E[\varepsilon_{ij}^4|\Theta, X] < C$  almost surely. (ii) At least one of the following holds:

- a  $\min_{ij} \text{Var}(\varepsilon_{ij}|\Theta, X) > c$ .
- b  $\varepsilon_{ij}$  can be decomposed as  $\varepsilon_{ij} = e_{ij}x_{ij}$  with  $\min_{ij} \text{Var}(e_{ij}|\Theta, X) > c$ .

Assumption 4.2(ii) is stated in a way that specifically covers the well-known matrix completion problem:

$$y_{ij}^* = \theta_{ij} + e_{ij},$$

where  $y_{ij}^*$  may not be observable, and  $x_{ij}$  indicates the observability for each element. Then  $\varepsilon_{ij} = e_{ij}x_{ij}$ .

The assumption below restricts the DGP of the design variable  $x_{ij}$ . The restrictions imposed are mild, and the assumption is stated so as to cover a variety of cases. Specifically, conditions (a)–(c) in Assumption 4.3 allow for various types of dependence among the  $x_{ij}$ .

ASSUMPTION 4.3 (DGP for  $x_{ij}$ ). (i)  $\max_{ij} |x_{ij}| < C$  and  $x_{ij}$  is independent of  $\Theta$ . (ii) At least one of the following holds:

- a  $x_{ij}^2$  does not vary across  $i \leq n$ .
- b  $x_{ij}^2$  is independent across  $(i, j)$ . In addition,  $E x_{ij}^2$  does not vary with  $i$ .
- c  $x_{ij} \in \{0, 1\}$ . Also, define  $\mathcal{B}_i := \{j \leq p : x_{ij} = 1\}$ . Then there is a set  $\bar{\mathcal{B}} \subseteq \{1, \dots, p\}$ , so that

(8)

$$\max_{i \leq n} \sum_{j=1}^p 1\{j \in \bar{\mathcal{B}} \Delta \mathcal{B}_i\} = o_P(d_{n,p}), \quad d_{n,p} := \left( \frac{\min\{n, p, \psi_{np}\}p}{(n+p)J + \|R\|_{(n)}^2} \right)^{J-(2+d+2b)},$$

where  $\bar{\mathcal{B}} \Delta \mathcal{B}_i = [\bar{\mathcal{B}} \cap \mathcal{B}_i^c] \cup [\bar{\mathcal{B}}^c \cap \mathcal{B}_i]$  is the symmetric difference of two sets, and  $d, b \geq 0$  are constants defined in Assumption 4.6 below.

Under Condition (ii)a, we can accommodate both conventional factor models by setting  $x_{ij} = 1$  for all  $i, j$  as well as conditional empirical factor models, where  $x_{ij} = x_j$ , with varying coefficients. An example of the latter is an asset pricing model with risk premia that vary across assets and over time where  $x_j$  represents the common time-varying market factor.

Condition (ii)b could cover examples of PCA with missing data under heterogeneous missing probabilities as in [Zhu, Wang and Samworth \(2022\)](#). In this case, we may take  $j$  to represent subjects and  $i$  to represent the index of repeated sampling within subject. The condition also accommodates scenarios where  $x_{ij}$  represents a treatment indicator where random assignment of subjects  $i$  to treatment states occurs independently in each period  $j$ . Such a structure may approximate some digital experimentation settings.

Condition (ii)c allows for some types of strong dependence in  $x_{ij}$  across both  $i$  and  $j$  but restricts  $x_{ij}$  to be binary as would be appropriate in missing data, matrix completion and treat-

ment assignment settings. In this condition, the set  $\mathcal{B}_i$  represents unit-specific “observation times” for unit  $i$ ; and the set  $\bar{\mathcal{B}}$  is common to all units. The quantity  $\max_{i \leq n} \sum_{j=1}^n 1\{j \in \bar{\mathcal{B}} \Delta \mathcal{B}_i\}$  thus measures the difference between the “unit specific” observation times and the “common” observation times. Condition (iii)c requires that these differences should be negligible. Hence, all units should be observed at approximately the same time. For instance, suppose every unit is observed most of the time in the sense that

$$\max_{i \leq n} \sum_{j=1}^p 1\{j : x_{ij} = 0\} = o_P(d_{n,p}).$$

Then Condition (ii)c holds with  $\bar{\mathcal{B}} = \{1, \dots, p\}$ .

Next, recall that  $v_j$  and  $u_i$  are respectively the  $j$ th right singular vector and the  $i$ th left singular vector of  $\Theta_0$ .

ASSUMPTION 4.4 (Incoherent singular vectors).

$$\mathbf{E} \max_{j \leq p} \|v_j\|^2 = O(Jp^{-1}), \quad \mathbf{E} \max_{i \leq n} \|u_i\|^2 = O(Jn^{-1}).$$

The incoherence condition ensures that information regarding the eigenspace accumulates as the dimension increases and allows us to apply our “rotation” argument to argue that estimating the eigenvector space is asymptotically unbiased. We provide low-level conditions that are sufficient for the incoherence condition in a treatment effects context where the low-rank matrix is formulated using nonparametric sieve representations in equation (17).

The next assumption places restrictions on various moments.

ASSUMPTION 4.5 (Moment bounds). There are matrices  $A_i, B_j$  whose eigenvalues are bounded away from zero and infinity, so that

$$\max_{i \leq n} \left\| \sum_{j=1}^p x_{ij}^2 v_j v_j' - A_i \right\| = o_P(J^{-1/2}), \quad \max_{j \leq p} \left\| \frac{n}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} x_{ij}^2 u_i u_i' - B_j \right\| = o_P(J^{-1/2}).$$

This should hold for  $\mathcal{S}$  being sets  $\{1, \dots, n\}, \mathcal{I}$  and  $\mathcal{I}^c$ .

Finally, we present the required conditions on  $\psi_{np}$ , the signal strength of the nonzero singular values. Recall that  $\psi_j(A)$  denotes the  $j$ th largest singular value of  $A$ . We allow the eigengap to change with  $J$ , depending on constants  $b, d \geq 0$ . This generality complicates statement of the condition but is needed to accommodate settings where the rank  $J$  is allowed to increase with sample sizes. We provide low-level conditions that are sufficient for the following assumption in the context of a treatment effect example in Lemma 5.2.

ASSUMPTION 4.6 (Signal-noise). There are constants  $b, d \geq 0$  such that:

(i)  $\psi_{np} \leq \psi_J(\Theta_0) < \psi_1(\Theta_0) \leq O_P(J^b \psi_{np})$  for a sequence  $\psi_{np} \rightarrow \infty$  that satisfies

$$n^{-1/2} p J^{7/2+2d+5b} + (p \vee n)^{3/4} J^{5/4+d+2b} = o(\psi_{np}).$$

(ii) Eigengap: There are  $c, C > 0$  and a sequence  $\psi_{np} \rightarrow \infty$  so that with probability approaching one,

$$\psi_j(\Theta_0) - \psi_{j+1}(\Theta_0) \geq c \psi_{np} J^{-d}, \quad j = 1, \dots, J.$$

(iii) The rank  $J$  satisfies

$$J^{3+2d+6b} = o_P(\min\{\sqrt{p}, \sqrt{n}, p/\sqrt{n}\}).$$

(iv) The low-rank approximation error matrix  $R = (r_{ij})_{n \times p}$  satisfies

$$\max_{ij} |r_{ij}|^2 (p \vee n)^2 J^{3+4b} = o(1).$$

**THEOREM 4.1.** *Suppose  $g$  is either dense or sparse, in the sense of (6) and (7). Suppose Assumptions 4.1–4.6 hold, and the nuclear-norm tuning parameter satisfies  $\nu > C(\sqrt{n+p})$  for some constant  $C > 0$ . Then for a fixed  $i \leq n$ ,*

$$\frac{\widehat{\theta}'_i g - \theta'_i g}{\sqrt{s_{np,1}^2 + s_{np,2}^2}} \rightarrow^d N(0, 1),$$

where, with  $L_j = \sum_{i=1}^n x_{ij}^2 \gamma_i \gamma'_i$  and  $\bar{B} = \sum_{j=1}^p (\mathbf{E} x_{ij}^2) v_j v'_j$ ,

$$s_{np,1}^2 := \sum_{j=1}^p \sum_{t=1}^n \text{Var}(\varepsilon_{tj} | \Theta, X) [\gamma'_t L_j^{-1} \gamma_t]^2 x_{tj}^2 g_j^2,$$

$$s_{np,2}^2 := \sum_{j=1}^p \text{Var}(\varepsilon_{ij} | \Theta, X) x_{ij}^2 [v'_j \bar{B}^{-1} V_0 g]^2.$$

To estimate the asymptotic variance, we need to preserve the rotation invariance property of the asymptotic variance. We therefore estimate  $s_{np,k}^2$  separately within subsamples and produce the final asymptotic variance estimator by averaging the results across subsamples. We consider the homoscedastic case where  $\text{Var}(\varepsilon_{ij} | \Theta, X) = \sigma_j^2$  for some constant  $\sigma_j^2$ ,  $j = 1, \dots, p$ . In this case, standard errors can be estimated as

$$\widehat{s}_{np,1}^2 := \frac{1}{4} \sum_{j=1}^p \sum_{t \notin \mathcal{I}} \widehat{\sigma}_j^2 [\widehat{\gamma}'_{t,\mathcal{I}} \widehat{L}_{j,\mathcal{I}}^{-1} \widehat{\gamma}_t]^2 x_{tj}^2 g_j^2 + \frac{1}{4} \sum_{j=1}^p \sum_{t \notin \mathcal{I}^c} \widehat{\sigma}_j^2 [\widehat{\gamma}'_{t,\mathcal{I}^c} \widehat{L}_{j,\mathcal{I}^c}^{-1} \widehat{\gamma}_t]^2 x_{tj}^2 g_j^2,$$

$$\widehat{s}_{np,2}^2 := \frac{1}{2} \sum_{j=1}^p \widehat{\sigma}_j^2 x_{ij}^2 [\widetilde{v}'_{j,\mathcal{I}} \widehat{B}_{\mathcal{I}}^{-1} \widetilde{V}_{\mathcal{I}} g]^2 + \frac{1}{2} \sum_{j=1}^p \widehat{\sigma}_j^2 x_{ij}^2 [\widetilde{v}'_{j,\mathcal{I}^c} \widehat{B}_{\mathcal{I}^c}^{-1} \widetilde{V}_{\mathcal{I}^c} g]^2,$$

$$\widehat{\sigma}_j^2 := \frac{1}{n} \sum_{t \notin \mathcal{I}} (y_{tj} - x_{tj} \cdot \widehat{\gamma}'_{t,\mathcal{I}} \widehat{v}_{j,\mathcal{I}})^2 + \frac{1}{n} \sum_{t \notin \mathcal{I}^c} (y_{tj} - x_{tj} \cdot \widehat{\gamma}'_{t,\mathcal{I}^c} \widehat{v}_{j,\mathcal{I}^c})^2,$$

where  $\widehat{L}_{j,\mathcal{I}} = \sum_{t \notin \mathcal{I}} x_{tj}^2 \widehat{\gamma}_t \widehat{\gamma}'_t$ , and  $\widehat{B}_{\mathcal{I}} = \sum_{j=1}^p x_{ij}^2 \widetilde{v}_{j,\mathcal{I}} \widetilde{v}'_{j,\mathcal{I}}$ , and  $\widehat{L}_{j,\mathcal{I}^c}$  and  $\widehat{B}_{\mathcal{I}^c}$  are defined similarly.

It is interesting to note that  $\sigma_{np}^2 := s_{np,1}^2 + s_{np,2}^2 = O_P(\frac{1}{n} \|g\|^2 + \frac{1}{p})$  in the case of fixed  $J$ . Thus, in this setting, the scaling of the asymptotic variance depends heavily on  $\|g\|^2$ .

**4.2. Semiparametric efficiency.** The semiparametric efficiency bound for the case of sparse  $g$  was established by [Chen et al. \(2019\)](#) (Lemma 2) in matrix completion settings and by [Iwakura and Okui \(2014\)](#) (Theorem 4.5) in pure factor models. Our asymptotic variance attains these previously established bounds if  $e_{ij}$  is i.i.d. homoscedastic Gaussian, so we do not further discuss semiparametric efficiency in the sparse setting.

We now provide a semiparametric efficiency bound in the case of dense  $g$  and verify that our estimator achieves this bound. For concreteness, suppose we are interested in  $h(\Theta) = \theta'_1 g$  where  $\theta'_1$  is the first row of  $\Theta$  and  $g$  is dense. In providing our result, we will allow for a wide range of distributions for  $x_{1j}$  while maintaining the assumption that the error term is Gaussian to make calculation tractable.

Specifically, we suppose that  $x_{ij}$  follows the distribution  $f$  and  $e_{ij} \sim N(0, \sigma^2)$  are independent across  $(i, j)$ . Let  $\mu_f = E_f x_{1j}^2$ . Under Assumptions 4.1–4.6, the dominant term in the asymptotic variance is

$$\begin{aligned} s_{np,2}^2 &= \sigma^2 \sum_{j=1}^p x_{ij}^2 [v'_j \bar{B}^{-1} V'_0 g]^2 \\ &= s_*^2(\Theta, f, \sigma) + o_P(s_{np,2}^2) \quad \text{where } s_*^2(\Theta, f, \sigma) = \sigma^2 \mu_f^{-1} \|V'_0 g\|^2, \end{aligned}$$

and we also have  $s_{np,1}^2 = o_P(s_{np,2}^2)$ . Hence,  $\hat{\theta}'_i g - \theta'_i g = O_P(\|V'_0 g\|)$  with asymptotic variance

$$s_{np,1}^2 + s_{np,2}^2 = s_*^2(\Theta, f, \sigma)(1 + o_P(1))$$

in this case.

The following result verifies that  $s_*^2(\Theta, f, \sigma)$  matches with the semiparametric efficiency bound. The notion of semiparametric efficiency in the presence of high-dimensional nuisance parameters is adopted from Janková and van de Geer (2018). The idea is to derive the asymptotic Cramér–Rao bound for asymptotically unbiased estimators, and needs to be formally established in the high-dimensional setting. Our result is novel relative to Janková and van de Geer (2018) because they deal with sparse models and our setting has low-rank matrices as the nuisance parameters.

**THEOREM 4.2.** *Consider  $h(\Theta) = \theta'_1 g$ , where  $\theta'_1$  is the first row of  $\Theta$  and  $g$  is dense. Let  $x_{ij} \sim f$  and  $e_{ij} \sim N(0, \sigma^2)$  be independent across  $(i, j)$ . Define*

$$\mathcal{M} = \{(A, f, \sigma) : \text{rank}(A) \leq J, \text{ Assumptions 4.1–4.6 hold}\}.$$

*Suppose that  $T(Y, X)$  is an asymptotically unbiased estimator of  $h(\Theta)$  in the sense that  $E_{(\Theta, f, \sigma)} T(Y, X) - h(\Theta) = o(s_*(\Theta, f, \sigma))$  where  $E_{(\Theta, f, \sigma)}$  denotes the expectation with respect to a given parameter  $(\Theta, f, \sigma)$ . Then for any sequence of  $(\Theta, f, \sigma) \in \mathcal{M}$ ,*

$$\liminf_{n, p \rightarrow \infty} \frac{E_{(\Theta, f, \sigma)} [T(Y, X) - h(\Theta)]^2}{s_*^2(\Theta, f, \sigma)} \geq 1.$$

**4.3. The role of spiked singular-values and incoherence.** Two key conditions that underlie our main results are the incoherence condition, Assumption 4.4 and the spiked singular-value (SSV) condition, Assumption 4.6. We demonstrate the role of these conditions by providing minimax theory for estimating  $\theta'_i g$  for a sparse or dense  $g$  without imposing SSV or incoherence, in a simple matrix completion problem where the missing indicators  $x_{ij}$  are independent Bernoulli random variables.

Define the following set of low-rank matrices:

$$\mathcal{S} = \left\{ A \in \mathbb{R}^{n \times p} : \text{rank}(A) \leq J \text{ and } \max_{1 \leq i \leq n} \max_{1 \leq j \leq p} |A_{ij}| \leq c_1 \right\}$$

for a constant  $c_1 > 0$  and for  $J \geq 1$ . Here,  $J$  is allowed to be either a fixed constant or a sequence tending to infinity.

We prove the following result for matrix completion over the space  $\mathcal{S}$ . Let  $y_{ij} = x_{ij} \theta_{ij} + e_{ij}$ , where  $x_{ij} \sim \text{Bernoulli}(\rho_j)$  and  $e_{ij} \sim N(0, \sigma_{ij}^2)$  are independent across  $(i, j)$ . Suppose that there are constants  $c_2, \dots, c_6 > 0$  such that  $\rho_j \in (c_2, 1 - c_2)$  and  $\sigma_{ij} \in (c_3, c_4)$  for any  $(i, j)$ . Let  $\rho = (\rho_1, \dots, \rho_p)'$  and  $\sigma = \{\sigma_{ij}\}_{1 \leq i \leq n, 1 \leq j \leq p}$ . In the theorem below,  $T$  represents any measurable function of the data, typically regarded as an “estimator” for  $h(\Theta) = \theta'_1 g$ .

**THEOREM 4.3 (Minimax rate).** *Consider estimating  $h(\Theta) = \theta'_1 g = \sum_{j=1}^p \theta_{1j} g_j$ , and let  $P_{(\Theta, f, \sigma)}$  denote the probability measure with respect to a given parameter  $(\Theta, f, \sigma)$ . We have the following results:*

1. *Sparse g: Let  $g_1 = 1$  and  $g_j = 0$  for  $j \geq 2$ , that is,  $h(\Theta) = \theta_{11}$ . Then*

$$(9) \quad \inf_T \sup_{\Theta \in \mathcal{S}} P_{(\Theta, \rho, \sigma)}(|T - h(\Theta)| > \kappa) > 1/4,$$

where  $\kappa > 0$  is a constant depending on  $(c_1, c_3)$  and  $\inf_T$  is taken over all measurable functions of the data  $(X, Y)$ .

2. *Dense g: Let  $|g_j| \in [c_5/p, c_6/p]$  for all  $j \in \{1, \dots, p\}$ . Then*

$$(10) \quad \inf_T \sup_{\Theta \in \mathcal{S}} P_{(\Theta, \rho, \sigma)}(|T - h(\Theta)| > \kappa p^{-1/2}) > 1/4,$$

where  $\kappa > 0$  is a constant depending on  $(c_1, c_3, c_5)$  and  $\inf_T$  is taken over all measurable functions of the data  $(X, Y)$ .

Theorem 4.3 gives the minimax rate *without* SSV and the incoherence condition. It provides a similar intuition to Koltchinskii, Löffler and Nickl (2020). For instance, (9) shows that it is impossible to guarantee entrywise consistency for sparse  $g$  in the considered setting without SSV or incoherence.

In addition, equation (10) implies that the rate  $O_p(p^{-1/2})$  is minimax optimal for estimating dense averages in the absence of SSV and incoherence. This rate of convergence is slower than that obtained in Theorem 4.1, which makes use of SSV and incoherence. For instance, in the factor model with a finite number of strong factors, Theorem 4.1 implies that the rate of convergence can be as fast as  $\frac{1}{p} \sum_j \widehat{\theta}_{ij} - \frac{1}{p} \sum_j \theta_{ij} = O_p(\frac{1}{\sqrt{np}} + \frac{1}{p})$ .<sup>2</sup>

These minimax results for estimating linear combinations of elements of a low-rank matrix without SSV and incoherence are new to the literature. The result closest to ours is Koltchinskii, Löffler and Nickl (2020), which provides minimax rates for estimating linear functionals of the eigenvectors of low-rank matrices. They show that the minimax optimal rate can be slow if the SSV condition does not hold. Other results on the minimax bounds for learning an eigenspace can be found in Berthet and Rigollet (2013), Birnbaum et al. (2013) and Cai, Ma and Wu (2013).

**4.4. Dense functional inference without SSV and incoherence.** When  $g$  is a vector of dense weights, the second minimax result in Theorem 4.3 suggests that consistency can be achieved without the SSV and incoherence conditions at the cost of a slower rate of convergence. For completeness, we introduce an alternative estimator that could be used when one does not wish to impose these assumptions.

Specifically, suppose  $g = (g_1, \dots, g_p)' \in \mathbb{R}^p$  is a vector of dense weights as defined in (7), and we are interested in the functional  $h_i(\Theta) := \theta'_i g$ . We propose the following estimator in the spirit of inverse probability weighting:

$$\widehat{h_i(\Theta)} = \sum_{j=1}^p \frac{g_j y_{ij} x_{ij}}{\widehat{\mu}_{j,i}^2}, \quad \widehat{\mu}_{j,i}^2 = \frac{1}{n-1} \sum_{k \neq i} x_{kj}^2.$$

---

<sup>2</sup>This rate holds if the factors have zero mean so that  $V'_0 g = \frac{1}{p} \sum_{j=1}^p v_j = O_p(p^{-1})$ , which is the case for no-intercept factor models. Strictly speaking, this setting was ruled out by Assumption 4.1, which requires  $\|V'_0 g\| \geq cp^{-1/2}$ . However, Assumption 4.1 is used only for obtaining the asymptotic distribution. This assumption can be relaxed when only the rate of convergence is of interest.

Note that this estimator does not require knowing the rank or even that the rank is consistently estimable. It is defined as the weighted average of the  $i$ th row of  $Y$  and  $X$  with weight proportional to a leave-one-out estimator of the inverse of  $\mu_j^2 := \text{Ex}_{ij}^2$ .

Let

$$W_{ij} := x_{ij}\varepsilon_{ij} + x_{ij}^2\theta_{ij}.$$

**THEOREM 4.4.** *Let  $g$  be dense in the sense of (7), and assume  $\text{Ex}_{ij}\varepsilon_{ij} = 0$ . Suppose  $W_{ij}$  is independent over  $j$  and that  $\text{EW}_{ij}^4 < C$ ,  $\text{Ex}_{ij}^2 > c > 0$ , and  $\text{Var}(W_{ij}) > c > 0$ . In addition, suppose  $\sqrt{p} \log p = o(n)$ . Then*

$$s_n^{-1} \sqrt{p} [\widehat{h}_i(\Theta) - \theta_i'g] \rightarrow^d N(0, 1),$$

where  $s_n^2 = p \sum_{j=1}^p g_j^2 (\text{Ex}_{ij}^2)^{-2} \text{Var}(W_{ij})$ .

**5. Application to heterogeneous treatment effects.** As an important illustration, we show how to apply our framework in a treatment effects setting. Suppose that, for each time  $j = 1, \dots, p$  and each unit  $i = 1, \dots, n$ , there is a pair of potential outcomes

$$(11) \quad Y_{ij}(m) = h_{j,m}(\eta_i) + e_{ij}(m), \quad m \in \{0, 1\}.$$

Here,  $m$  denotes treatment ( $m = 1$ ) or control ( $m = 0$ ) state. In any time period  $j$  and for any unit  $i$ , we observe either  $Y_{ij}(1)$  or  $Y_{ij}(0)$ , but not both, depending on the unit’s realized treatment state in that period. The treatment effect depends on time-varying functions  $h_{j,m}(\cdot)$  of unit specific state variable  $\eta_i$ ; both  $h_{j,m}(\cdot)$  and  $\eta_i$  may be unobservable and random. For clarity, we focus on the scenario where the goal is to perform statistical inference on a long-run treatment effect for a given unit  $i$ :

$$\tau_i := \frac{1}{p} \sum_{j=1}^p v_{ij},$$

where  $v_{ij} = h_{j,1}(\eta_i) - h_{j,0}(\eta_i)$  is the treatment effect for unit  $i$  at time  $j$ .

Define the treatment status indicator

$$x_{ij}(m) = 1\{\text{unit } i \text{ at period } j \text{ is in state } m\} = 1\{Y_{ij}(m) \text{ is observable}\}.$$

Consider the following treatment scenario. Suppose the entire time span  $\{1, 2, \dots, p\}$  is divided into two periods,

$$T_0 = \{1, \dots, p_0\} \quad \text{and} \quad T_1 = \{p_0 + 1, \dots, p\},$$

where both  $p_0$  and  $p_1 := p - p_0$  are large and both periods are known. We assume

$$(12) \quad \begin{aligned} \max_{i \leq n} \sum 1\{j \in T_0 : x_{ij}(0) = 0\} &= o_P(d_{n,p_0}), \\ \max_{i \leq n} \sum 1\{j \in T_1 : x_{ij}(1) = 0\} &= o_P(d_{n,p_1}), \end{aligned}$$

where  $d_{n,p_0}$  and  $d_{n,p_1}$  are slowly growing sequences defined in (8). That is, each unit is in the control state during most periods in  $T_0$ , and each unit is the treatment state during most periods in  $T_1$ . We thus refer to  $T_0$  and  $T_1$ , respectively, as the “control period” and the “treatment period.” We refer to this treatment scenario as “systematic treatment,” and note that treatment assignments are strongly dependent in this setting, which results in an important difference from much of the literature on inference in matrix completion settings. In terms of our formal conditions, this scenario corresponds to the case of Assumption 4.1(ii)c.

5.1. *Treatment effect inference.* Let  $\theta_{ij}(m) := h_{j,m}(\eta_i)$ . We can then rewrite the model for potential outcomes (11) as

$$(13) \quad y_{ij}(0) = \theta_{ij}(0)x_{ij}(0) + \varepsilon_{ij}(0), \quad j \in T_0,$$

$$(14) \quad y_{ij}(1) = \theta_{ij}(1)x_{ij}(1) + \varepsilon_{ij}(1), \quad j \in T_1,$$

where  $y_{ij}(m) = Y_{ij}(m)x_{ij}(m)$  and  $\varepsilon_{ij}(m) = e_{ij}(m)x_{ij}(m)$ . Let  $\Theta(m)$  denote the  $n \times p$  matrix of  $(\theta_{ij}(m))_{n \times p}$ . As, for example, previously note by [Athey et al. \(2021\)](#), it is then clear that recovering elements of  $\Theta(m)$  is equivalent to solving a matrix completion problem.

In Section 5.2, we provide sufficient conditions to establish that  $\Theta(m)$  is an approximate low-rank matrix that satisfies both the SSV and incoherence conditions. Under these conditions, we can then estimate treatment effects by simply applying Algorithm 2.1 twice—once using the data from period  $T_0$  and once using the data from period  $T_1$ .

*Step 1:* Apply Algorithm 2.1 to (13) to estimate  $\Theta(0)$ .

*Step 2:* Apply Algorithm 2.1 to (14) to estimate  $\Theta(1)$ .

*Step 3:* Make inference on the treatment effects from the estimated  $\Theta(1) - \Theta(0)$ .

Let  $\hat{\theta}_{ij}(m)$  denote the  $(i, j)$  element of the estimated matrix  $\Theta(m)$ . The average treatment effect estimator is then given by

$$\hat{\tau}_i := \frac{1}{p_1} \sum_{j \in T_1} \hat{\theta}_{ij}(1) - \frac{1}{p_0} \sum_{j \in T_0} \hat{\theta}_{ij}(0).$$

It is straightforward to extend Theorem 4.1 to this context, which leads to the asymptotic distribution of the estimated treatment effects. Formal results are to be presented in Section 5.4.

5.2. *The low-rank approximation.* We show that the matrix formed from elements  $h_{j,m}(\eta_i)$  can be approximated by a low-rank matrix with a slowly growing rank. To aid in focusing on the main idea, we suppress the notation “ $m$ ” throughout this section.

Consider a family of time-varying functions  $h_j(\cdot)$  of subject-specific latent variables  $\eta_i$ . Let  $\Theta$  be the  $n \times p$  matrix obtained by setting the  $(i, j)$  element of  $\Theta$  to  $h_j(\eta_i)$ . Suppose  $h_j(\cdot)$  has a sieve approximation:

$$(15) \quad h_j(\eta_i) = \sum_{k=1}^J \lambda_{j,k} \phi_k(\eta_i) + r_{ij} = \lambda'_j \Phi_i + r_{ij},$$

where  $\Phi_i := (\phi_1(\eta_i), \dots, \phi_J(\eta_i))' \in \mathbb{R}^J$  is a set of sieve transformations of  $\eta_i$  using  $\phi_k(\cdot)$  as the basis functions,  $\lambda_j = (\lambda_{j,1}, \dots, \lambda_{j,J})'$  is the vector of sieve coefficients for  $h_j(\cdot)$ , and  $r_{ij}$  is the sieve approximation error. Write  $\Phi$  as the  $n \times J$  matrix of  $\Phi_i$ ,  $\Lambda$  as the  $p \times J$  matrix of  $\lambda_j$  and  $R$  as the  $n \times p$  matrix of  $r_{ij}$ . Then the matrix form of (15) is

$$\Theta = \underbrace{\Phi \Lambda'}_{\Theta_0} + R.$$

Clearly,  $\text{rank}(\Theta_0) \leq J$ , and there is a rotation matrix  $H$  so that columns of  $\Lambda H$  are the right singular vectors of  $\Theta_0$ . The error  $R$  is naturally present as the sieve approximation error, which will decrease as more elements are considered in the sieve approximation. It is then natural to consider sequences where  $J$  increases slowly with  $(n, p)$ .

We now illustrate how both the SSV and incoherence conditions can hold in this setting under sensible conditions on the functional space and the sieve bases. Suppose  $h_j$  belongs to a Hölder class: For some  $C, \beta, \alpha > 0$ ,

$$\left\{ h : \max_{\gamma_1 + \dots + \gamma_d = \beta} \left| \frac{\partial^\beta h(x)}{\partial x_1^{\gamma_1} \dots \partial x_d^{\gamma_d}} - \frac{\partial^\beta h(y)}{\partial y_1^{\gamma_1} \dots \partial y_d^{\gamma_d}} \right| \leq C \|x - y\|^\alpha, \text{ for all } x, y \right\}.$$

Further suppose that a common basis, such as polynomials or B-splines are considered. We will then have

$$\max_{ij} |r_{ij}| \leq C J^{-a}, \quad a = (\beta + \alpha) / \dim(\eta_i),$$

which can be made arbitrarily small for sufficiently smooth functions even if  $J$  grows slowly.

Now, suppose there exists a  $b \geq 0$  such that  $\psi_J(\Theta) \leq \psi_1(\Theta) \leq C J^b \psi_J(\Theta)$  for some  $C > 1$ . It is then easy to show that the sequence  $\psi_{np}$  can be taken as

$$\psi_{np} \asymp \sqrt{J^{-(2b+1)} \sum_{i=1}^n \sum_{j=1}^p h_j(\eta_i)^2}.$$

We then have that the top  $J$  singular values grow at this rate, which leads to the SSV condition.

Finally, write  $S_\Lambda = \frac{1}{p} \Lambda' \Lambda$ ,  $S_\Phi = \frac{1}{n} \Phi' \Phi$  and  $A = S_\Phi^{1/2} S_\Lambda S_\Phi^{1/2}$ . Also, let  $G_\Phi$  be a  $J \times J$  matrix whose columns are the eigenvectors of  $A$ , and let  $T$  be the diagonal matrix of corresponding eigenvalues. Letting  $H_\Phi := S_\Phi^{-1/2} G_\Phi$ , it can be verified that

$$\Theta_0 \Theta_0' \Phi H_\Phi = pn \Phi H_\Phi T \quad \text{and} \quad \frac{1}{n} (\Phi H_\Phi)' \Phi H_\Phi = I.$$

Thus, the columns of  $\frac{1}{\sqrt{n}} \Phi H_\Phi$  are the left singular vectors of  $\Theta_0$ , and the eigenvalues of  $npA$  equal the first  $J$  eigenvalues of  $\Theta_0' \Theta_0$ . Similarly, we can define  $H_\Lambda = S_\Lambda^{-1/2} G_\Lambda$  where  $G_\Lambda$  is a  $J \times J$  matrix whose columns are the eigenvectors of  $S_\Lambda^{1/2} S_\Phi S_\Lambda^{1/2}$ . Hence, we have

$$(16) \quad U_0 = n^{-1/2} \Phi H_\Phi, \quad V_0 = p^{-1/2} \Lambda H_\Lambda.$$

Thus,

$$(17) \quad \begin{aligned} \max_{i \leq n} \|u_i\| &\leq n^{-1/2} \max_{i \leq n} \|\Phi_i\| \psi_{\min}^{-1/2}(S_\Phi), \\ \max_{j \leq p} \|v_j\| &\leq p^{-1/2} \max_{j \leq p} \|\lambda_j\| \psi_{\min}^{-1/2}(S_\Lambda). \end{aligned}$$

It then follows that the incoherence condition holds as long as we can obtain proper upper bounds for  $\max_{j \leq p} \|\lambda_j\|$  and  $\max_{i \leq n} \|\Phi_i\|$ . For example, if  $\{h_j(\cdot) : j \leq p\}$  is further restricted to a Hilbert space with a uniform  $L_2$ - bound,

$$\max_{j \leq p} \sum_{k=1}^{\infty} \lambda_{j,k}^2 < \infty,$$

then  $\max_{j \leq p} \|\lambda_j\| < C$ .

We formalize the preceding discussion in the following assumption and lemma.

ASSUMPTION 5.1. (i)  $\max_{j \leq J} \sup_{\eta} |\phi_j(\eta)| < C$ ,  $E \psi_{\min}^{-1}(S_\Phi) < C$ , and  $\psi_{\min}^{-1}(S_\Lambda) < C$ .  
 (ii) The sieve approximation satisfies

$$\max_{ij} |r_{ij}| \leq C J^{-a}$$

for some  $a > 0$ .

(iii)  $\{h_j(\cdot) : j \leq p\}$  belong to ball  $\mathcal{H}(\mathcal{U}, \|\cdot\|_{L_2}, C)$  inside a Hilbert space spanned by the basis  $\{\phi_k : k = 1, \dots\}$  with a uniform  $L_2$ -bound  $C$ :

$$\sup_{h \in \mathcal{H}(\mathcal{U}, \|\cdot\|_{L_2})} \|h\| \leq C,$$

where  $\mathcal{U}$  is the support of  $\eta_i$ .



LEMMA 5.1. *Suppose Assumption 5.1 holds. Then:*

(i) *The minimum nonzero singular value  $\psi_{np}$  for  $\Theta_0 = \Phi \Lambda'$  can be taken as*

$$\psi_{np}^2 \asymp J^{-(2b+1)} \sum_{i=1}^n \sum_{j=1}^p h_j(\eta_i)^2, \quad m = 0, 1,$$

*which means  $\psi_J(\Theta_0) \geq c\psi_{np}$  for this choice of  $\psi_{np}$ .*

(ii) *The incoherence Assumption 4.4 holds.*

(iii) *The low-rank approximation error satisfies  $\|R\|_{(n)} \leq C(p \vee n)^{3/2} J^{-a}$ .*

5.3. *Reproducing kernel representation.* We now verify the eigengap condition: Let  $A = \frac{1}{pn} \Theta \Theta'$ . There are constants  $b, d \geq 0$  such that

$$(18) \quad \begin{aligned} \psi_1(A) / \psi_J(A) &\leq O_P(J^b), \\ \min_{k=1, \dots, J-1} \psi_k(A) - \psi_{k+1}(A) &\geq cJ^{-d}. \end{aligned}$$

Below we verify the above conditions when the treatment functions are generated from a Gaussian process.<sup>3</sup>

Suppose  $\eta_i$  are uniformly generated from  $[0, 1]$ , and functions  $h_j(\cdot)$  are independently generated from a Gaussian process with covariance kernel

$$K(\eta_1, \eta_2) = \text{Cov}(h_j(\eta_1), h_j(\eta_2)),$$

where  $K(\cdot, \cdot)$  is a continuous positive semidefinite kernel function supported on a compact set. In addition, suppose the associated integral operator

$$(Tf)(\cdot) = \int K(\cdot, \eta) f(\eta) d\eta$$

is positive semidefinite. Let  $\{\bar{\phi}_k(\cdot)\}$  and  $v_k \geq 0$  be the eigenfunctions and eigenvalues of  $T$ . Then by Mercer's theorem,  $\{\bar{\phi}_k(\cdot)\}$  is an orthonormal basis so that  $K$  has the following representation:

$$K(\eta_1, \eta_2) = \sum_{k=1}^{\infty} v_k \bar{\phi}_k(\eta_1) \bar{\phi}_k(\eta_2),$$

where the infinite sum can be approximated arbitrarily well by finite truncation  $J$  as  $J \rightarrow \infty$ .

Now consider the  $n \times n$  matrix  $\frac{1}{p} \Theta \Theta'$ , whose  $(i, l)$  element is

$$\frac{1}{p} \sum_{j=1}^p h_j(\eta_i) h_j(\eta_l) = K(\eta_i, \eta_l) + o_P(1) = \bar{\Phi}'_i D_\lambda \bar{\Phi}_l + o_P(1),$$

where  $\bar{\Phi}'_i = (\bar{\phi}_1(\eta_i), \dots, \bar{\phi}_J(\eta_i))$  and  $D_\lambda$  is a diagonal matrix of  $(v_1, \dots, v_J)$ . Also, because the  $h_j$  are independently generated from the Gaussian process, the  $o_P(1)$  terms are uniform over all elements. Thus, we have an approximate low-rank representation of  $\Theta \Theta'$ :

$$\Theta \Theta' = \left[ \sum_j h_j(\eta_i) h_j(\eta_l) \right]_{n \times n} \approx p \bar{\Phi} D_\lambda \bar{\Phi}'.$$

Because the columns of  $\bar{\Phi}$  are formed from eigenfunctions, its columns are approximately orthonormal bases as eigenvectors of  $\Theta \Theta'$ . Hence, the diagonals of  $D_\lambda$  are also approximately

<sup>3</sup>We are thankful to one of the referees for suggesting this case.

the top  $J$  eigenvalues of  $\frac{1}{np} \Theta \Theta'$ . This observation heuristically shows that the top eigenvalues of  $\frac{1}{np} \Theta \Theta'$  are approximately the same as those of the integral operator  $T$  associated with the reproducing kernel function.

Rigorously, we can verify this condition as follows. The conditions of the lemma below are required to hold for both  $m \in \{0, 1\}$  in our treatment effect setting.

LEMMA 5.2. *Suppose the eigenvalues of the integral operator  $T$  satisfy*

$$v_k = Mk^{-\alpha}, \quad k = 1, 2, \dots$$

for some  $M, \alpha > 0$ . Further, suppose  $\sqrt{\frac{\log n}{p}} + r_J + \frac{J}{\sqrt{n}} = o_P(J^{-\alpha-1})$ , where  $r_J := \sup_{\eta_1, \eta_2} |\sum_{k>J} v_k \bar{\phi}_k(\eta_1) \bar{\phi}_k(\eta_2)|$ . Then the eigengap condition (18) holds. Specifically, let  $A = \frac{1}{pn} \Theta \Theta'$ ,

$$(19) \quad \begin{aligned} \psi_1(A) / \psi_J(A_m) &\leq O_P(J^\alpha), \\ \min_{k=1, \dots, J-1} \psi_k(A_m) - \psi_{k+1}(A) &\geq cJ^{-(\alpha+1)}. \end{aligned}$$

5.4. *Inference for treatment effects under systematic assignment.* Building on the previous subsections, suppose  $h_{j,m}(\eta_i)$  has the following sieve representation:

$$h_{j,m}(\eta_i) = \sum_{k=1}^J \lambda_{j,k,m} \phi_k(\eta_i) + r_{ij}(m), \quad m = 0, 1.$$

We then have that the matrix  $\Theta(m) := (\theta_{ij}(m))_{n \times p_m}$  admits an approximate low-rank structure for each  $m \in \{0, 1\}$ :

$$(20) \quad \Theta(m) = \Theta_0(m) + R(m), \quad \Theta_0(m) = \Phi \Lambda'_m, R(m) = (r_{ij}(m))_{n \times p_m},$$

where  $\Lambda_m$  is the  $p \times J$  matrix of  $\lambda_{j,k,m}$ .

Note that  $\hat{\tau}_i$  estimating a sensible average treatment effect relies on an additional stability assumption. Define for  $m \in \{0, 1\}$ ,

$$\zeta_{ij}(m) := x_{ij}(m) v_j(m)' \bar{B}(m)^{-1} \frac{1}{p_m} \sum_{j \in T_m} v_j(m),$$

where  $\bar{B}(m) = \sum_{j \in T_m} x_{ij}(m) v_j(m) v_j(m)'$ . Applying the analysis of Theorem 4.1, we have

$$\begin{aligned} \hat{\tau}_i - \tau_i &= \sum_{j \in T_1} e_{ij} \zeta_{ij}(1) - \sum_{j \in T_0} e_{ij} \zeta_{ij}(0) + o_P(\min\{p_0, p_1\}^{-1/2}) \\ &\quad + \left( \frac{1}{p_1} \sum_{j \in T_1} \theta_{ij}(1) - \frac{1}{p} \sum_{j=1}^p \theta_{ij}(1) \right) - \left( \frac{1}{p_0} \sum_{j \in T_0} \theta_{ij}(0) - \frac{1}{p} \sum_{j=1}^p \theta_{ij}(0) \right). \end{aligned}$$

This expansion yields the asymptotic distribution of  $\hat{\tau}_i$  under the condition that the second line on the right-hand side is bounded by  $o_P(\min\{p_0, p_1\}^{-1/2})$ . That is, we need stability of treatment and control averages in the sense that the average of  $\theta_{ij}(0)$  and  $\theta_{ij}(1)$  obtained over the respective subsamples does not deviate too far from the infeasible average that would be obtained looking over the entire sample period.

**THEOREM 5.3.** *Suppose Assumptions 4.1, 4.5, 4.6 hold. Suppose Assumption 5.1 holds for  $h_{j,0}$  and  $h_{j,1}$ . In addition, suppose*

$$\begin{aligned} \frac{1}{p_1} \sum_{j \in T_1} \theta_{ij}(1) - \frac{1}{p} \sum_{j=1}^p \theta_{ij}(1) &= o_P(\min\{p_0, p_1\}^{-1/2}) \quad \text{and} \\ \frac{1}{p_0} \sum_{j \in T_0} \theta_{ij}(0) - \frac{1}{p} \sum_{j=1}^p \theta_{ij}(0) &= o_P(\min\{p_0, p_1\}^{-1/2}). \end{aligned}$$

Let

$$\bar{s}_{np,i}^2 := \sum_{j \in T_0} \text{Var}(e_{ij}|X, \eta) \zeta_{ij}(0)^2 + \sum_{j \in T_1} \text{Var}(e_{ij}|X, \eta) \zeta_{ij}(1)^2.$$

Suppose there is a constant  $c > 0$  so that  $\bar{s}_{np,i}^2 \min\{p_0, p_1\} > c$  with probability approaching one. Then as  $n, p_0, p_1 \rightarrow \infty$ ,

$$\frac{\hat{\tau}_i - \tau_i}{\bar{s}_{np,i}} \rightarrow^d N(0, 1).$$

**6. Simulations.** We now illustrate the performance of our inferential approach through a small simulation study in the systematic treatment assignment setting. We report results for  $n = p = 400$ .

To generate data, we first divide the period of observation  $\{1, \dots, p\}$  equally into two periods  $T_0$  and  $T_1$  each consisting of  $p_m = p/2$  observation times. To generate  $x_{ij}(m)$ , we generate  $n_i$  integers  $j_1, \dots, j_{n_i}$  without replacement to form a set  $A_i(m) = \{j_1, \dots, j_{n_i}\} \subset T_m$ . The number  $n_i \leq N_0$  is uniformly generated to be less than a predetermined number  $N_0 \in \{p_m^{1/2}, p_m^{1/3}, p_m^{1/4}\}$ . We then set

$$x_{ij}(m) = \begin{cases} 0 & \text{if } j \in A_i(m), \\ 1 & \text{if } j \notin A_i(m). \end{cases}$$

Hence, for each unit  $i$ ,  $x_{ij}(m) = 1$ , with up to  $N_0$  exceptions, throughout period  $T_m$  whose total length is  $p_m$ . In addition, we generate the noise  $\varepsilon_{ij}$  independently across both  $(i, j)$  and  $\varepsilon_{ij}(m) \sim \mathcal{N}(0, \sigma_\varepsilon^2)$  for  $\sigma_\varepsilon = 1$ .

One of the key conditions in this scenario is that the treatment effect should be stable in the sense that  $\frac{1}{p} \sum_{j=1}^p \theta_{ij}(m)$  can be well approximated by  $\frac{1}{p_m} \sum_{j \in T_m} \theta_{ij}(m)$ . We thus consider the simplest possible setting where this condition holds by generating time invariant treatment functions:

$$h_0(\eta_i) = \sum_{k=1}^{\infty} \frac{|W_k|}{k^a} \sin(k\eta_i), \quad h_1(\eta_i) = \sum_{k=1}^{\infty} \frac{(|W_k| + 2)}{k^a} \sin(k\eta_i).$$

Here,  $\eta_i \sim \text{Uniform}[-1, 1]$ ,  $W_k \sim \mathcal{N}(0, 1)$  and the noise is  $e_{ij} \sim \mathcal{N}(0, 1)$ . The power parameter  $a > 1$  quantifies the decay speed of the sieve coefficients. The left panel of Figure 1 plots both functions. To illustrate the heterogeneity of  $\tau_i$  across  $i$ , Figure 1 also plots the histogram of  $\tau_i = h_1(\eta_i) - h_0(\eta_i)$  for  $i = 1, \dots, n$  for one simulation replication.

In terms of implementation of our procedure, we also need  $J$  and  $\nu$ . We do not attempt to infer the rank  $J$  from the data. Rather, we look at estimates based on four prespecified values of the rank:  $J = 1, \dots, 4$ . We set the parameter  $\nu$  for the nuclear-norm penalized optimization through a simple plug-in procedure. Specifically, we set

$$(21) \quad \nu = 2.2 \bar{Q}(\|Z \circ X(m)\|; 0.95),$$

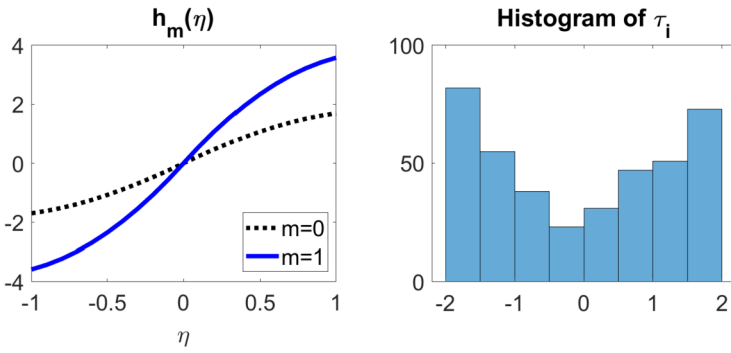


FIG. 1. Treatment functions  $h_m(\eta)$  with  $a = 3$  and histogram of  $\tau_i$  over  $i = 1, \dots, n$ .

where  $\bar{Q}(W; q)$  denotes the  $q$ th quantile of a random variable  $W$  and  $Z$  is an  $n \times p_m$  matrix whose elements  $z_{ij}$  are generated as  $\mathcal{N}(0, \hat{\sigma}_e^2)$  independent across  $(i, j)$  for some estimated  $\hat{\sigma}_e^2$ .<sup>4</sup> This choice can be motivated as in Belloni and Chernozhukov (2013) and Chernozhukov et al. (2018).

We report simulation coverage probabilities of 95% confidence intervals for  $\tau_1$  formed using estimated standard errors based on 1000 simulation replications in Table 1. Overall, the derived asymptotic distributions seem to provide reasonable approximations to the finite-sample distributions under our simulation settings, and the good performance appears quite robust to the choice of  $J$  in this simulation.

**7. Discussion.** In the treatment effect study, we approach the problem by separately applying our generic approach to an initial period where most units are in the control state and a subsequent period where most units are in the treatment state. The illustration is by no

TABLE 1  
Systematic Assignments. Coverage Probabilities of the treatment effect  $\tau_i$

$N_0$	Power $a$	$J$			
		1	2	3	4
$p_m^{1/2}$	4	0.952	0.950	0.949	0.948
	3	0.947	0.943	0.943	0.942
	2	0.952	0.950	0.948	0.949
$p_m^{1/3}$	4	0.950	0.950	0.947	0.945
	3	0.948	0.946	0.945	0.943
	2	0.952	0.950	0.948	0.947
$p_m^{1/4}$	4	0.954	0.952	0.949	0.946
	3	0.954	0.952	0.951	0.950
	2	0.955	0.956	0.950	0.951

Note: This table reports the simulated coverage probability of 95% confidence intervals. The rank  $J$  equals the sieve dimension used. Power  $a$  quantifies the decay rate of the sieve coefficients  $\lambda_{j,k} \sim k^{-a}$ . Finally,  $N_0$  controls the number of “exceptions” over time (the maximum number of treated during “control period” and the maximum number of controlled during “treatment period.”).

<sup>4</sup>We set  $\hat{\sigma}_e^2$  by obtaining an initial guess,  $\tilde{\sigma}_e^2$ , from estimating the simple model  $y_{ij} = x_{ij}\theta_i + \sigma_e^{-1}u_{ij}$  where  $\text{Var}(u_{ij}) = 1$ . We then obtain an initial solution to the nuclear-norm regularized optimization problem with tuning parameter set as in (21) with  $z_{ij} \sim N(0, \tilde{\sigma}_e^2)$ . Letting  $\tilde{\theta}_{ij}$  denote the nuclear-norm regularized estimator obtained with this initial tuning. We then set  $\hat{\sigma}_e^2 = \frac{1}{np} \sum_{ij} \tilde{\varepsilon}_{ij}^2$ , where  $\tilde{\varepsilon}_{ij} = y_{ij} - x_{ij}\tilde{\theta}_{ij}$ .

means exhaustive. For example, one could consider imposing invariance such that the low-rank matrices in the treatment and control states share the same singular space. Under this structure, estimation could be coupled rather than treating estimation as two separate problems. The problem could also be viewed as a tensor completion problem where the third dimension is the treatment assignment, so that  $\theta_{ij}(m)$  denotes the  $(i, j, m)$  element of the low-rank tensor  $\Theta = (\theta_{ij}(m))_{n \times p \times 2}$ . Then we may extend the proposed inference procedure to the tensor-recovery setting. Finally, Chen et al. (2019) used an “auxiliary leave-one-out” argument without actual sample splitting in the matrix completion setting with homogeneous missing data. It may be worth extending their approach to the current context. We expect each of these directions would be interesting for future research.

## SUPPLEMENTARY MATERIAL

**Inference for low-rank models: Online appendix** (DOI: [10.1214/23-AOS2293SUPP.pdf](https://doi.org/10.1214/23-AOS2293SUPP.pdf)). The Online Appendix Chernozhukov et al. (2023) contains all proofs.

## REFERENCES

- ABBE, E., FAN, J., WANG, K. and ZHONG, Y. (2020). Entrywise eigenvector analysis of random matrices with low expected rank. *Ann. Statist.* **48** 1452–1474. MR4124330 <https://doi.org/10.1214/19-AOS1854>
- ATHEY, S., BAYATI, M., DOUDCHENKO, N., IMBENS, G. and KHOSRAVI, K. (2021). Matrix completion methods for causal panel data models. *J. Amer. Statist. Assoc.* **116** 1716–1730. MR4353709 <https://doi.org/10.1080/01621459.2021.1891924>
- BELLONI, A. and CHERNOZHUKOV, V. (2013). Least squares after model selection in high-dimensional sparse models. *Bernoulli* **19** 521–547. MR3037163 <https://doi.org/10.3150/11-BEJ410>
- BERTHET, Q. and RIGOLLET, P. (2013). Optimal detection of sparse principal components in high dimension. *Ann. Statist.* **41** 1780–1815. MR3127849 <https://doi.org/10.1214/13-AOS1127>
- BIRNBAUM, A., JOHNSTONE, I. M., NADLER, B. and PAUL, D. (2013). Minimax bounds for sparse PCA with noisy high-dimensional data. *Ann. Statist.* **41** 1055–1084. MR3113803 <https://doi.org/10.1214/12-AOS1014>
- CAI, T. T., MA, Z. and WU, Y. (2013). Sparse PCA: Optimal rates and adaptive estimation. *Ann. Statist.* **41** 3074–3110. MR3161458 <https://doi.org/10.1214/13-AOS1178>
- CANDÈS, E. J. and PLAN, Y. (2010). Matrix completion with noise. *Proc. IEEE* **98** 925–936.
- CANDÈS, E. J. and RECHT, B. (2009). Exact matrix completion via convex optimization. *Found. Comput. Math.* **9** 717–772. MR2565240 <https://doi.org/10.1007/s10208-009-9045-5>
- CHEN, Y., CHI, Y., FAN, J., MA, C. and YAN, Y. (2020). Noisy matrix completion: Understanding statistical guarantees for convex relaxation via nonconvex optimization. *SIAM J. Optim.* **30** 3098–3121. MR4167625 <https://doi.org/10.1137/19M1290000>
- CHEN, Y., FAN, J., MA, C. and YAN, Y. (2019). Inference and uncertainty quantification for noisy matrix completion. *Proc. Natl. Acad. Sci. USA* **116** 22931–22937. MR4036123 <https://doi.org/10.1073/pnas.1910053116>
- CHERNOZHUKOV, V., HANSEN, C., LIAO, Y. and ZHU, Y. (2018). Inference for heterogeneous effects using low-rank estimations. ArXiv preprint. Available at arXiv:1812.08089.
- CHERNOZHUKOV, V., HANSEN, C., LIAO, Y. and ZHU, Y. (2023). Supplement to “Inference for low-rank models.” <https://doi.org/10.1214/23-AOS2293SUPP>
- CHERNOZHUKOV, V., HANSEN, C. and SPINDLER, M. (2015). Valid post-selection and post-regularization inference: An elementary, general approach. *Ann. Rev. Econ.* **7** 649–688.
- DRAY, S. and JOSSE, J. (2015). Principal component analysis with missing values: A comparative survey of methods. *Plant Ecol.* **216** 657–667.
- FAN, J., GUO, J. and ZHENG, S. (2022). Estimating number of factors by adjusted eigenvalues thresholding. *J. Amer. Statist. Assoc.* **117** 852–861. MR4436317 <https://doi.org/10.1080/01621459.2020.1825448>
- GROSS, D. (2011). Recovering low-rank matrices from few coefficients in any basis. *IEEE Trans. Inf. Theory* **57** 1548–1566. MR2815834 <https://doi.org/10.1109/TIT.2011.2104999>
- HASTIE, T., MAZUMDER, R., LEE, J. D. and ZADEH, R. (2015). Matrix completion and low-rank SVD via fast alternating least squares. *J. Mach. Learn. Res.* **16** 3367–3402. MR3450542
- IWAKURA, H. and OKUI, R. (2014). Asymptotic efficiency in factor models and dynamic panel data models. Available at SSRN 2395722.
- JAIN, P., NETRAPALLI, P. and SANGHAVI, S. (2013). Low-rank matrix completion using alternating minimization (extended abstract). In *STOC’13—Proceedings of the 2013 ACM Symposium on Theory of Computing* 665–674. ACM, New York. MR3210828 <https://doi.org/10.1145/2488608.2488693>

- JANKOVÁ, J. and VAN DE GEER, S. (2018). Semiparametric efficiency bounds for high-dimensional models. *Ann. Statist.* **46** 2336–2359. MR3845020 <https://doi.org/10.1214/17-AOS1622>
- JANKOVÁ, J. and VAN DE GEER, S. (2021). De-biased sparse PCA: Inference for eigenstructure of large covariance matrices. *IEEE Trans. Inf. Theory* **67** 2507–2527. MR4282369 <https://doi.org/10.1109/TIT.2021.3059765>
- KESHAVAN, R. H., MONTANARI, A. and OH, S. (2010). Matrix completion from a few entries. *IEEE Trans. Inf. Theory* **56** 2980–2998. MR2683452 <https://doi.org/10.1109/TIT.2010.2046205>
- KLOPP, O. (2014). Noisy low-rank matrix completion with general sampling distribution. *Bernoulli* **20** 282–303. MR3160583 <https://doi.org/10.3150/12-BEJ486>
- KOLTCHINSKII, V., LÖFFLER, M. and NICKL, R. (2020). Efficient estimation of linear functionals of principal components. *Ann. Statist.* **48** 464–490. MR4065170 <https://doi.org/10.1214/19-AOS1816>
- KOLTCHINSKII, V., LOUNICI, K. and TSYBAKOV, A. B. (2011). Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *Ann. Statist.* **39** 2302–2329. MR2906869 <https://doi.org/10.1214/11-AOS894>
- NEGAHBAN, S. and WAINWRIGHT, M. J. (2011). Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *Ann. Statist.* **39** 1069–1097. MR2816348 <https://doi.org/10.1214/10-AOS850>
- NEYMAN, J. (1959). Optimal asymptotic tests of composite statistical hypotheses. In *Probability and Statistics: The Harald Cramér Volume* (U. Grenander, ed.) 213–234. Almqvist & Wiksell, Stockholm. MR0112201
- ONATSKI, A. (2010). Determining the number of factors from empirical distribution of eigenvalues. *Rev. Econ. Stat.* **92** 1004–1016.
- RECHT, B. (2011). A simpler approach to matrix completion. *J. Mach. Learn. Res.* **12** 3413–3430. MR2877360
- ROHDE, A. and TSYBAKOV, A. B. (2011). Estimation of high-dimensional low-rank matrices. *Ann. Statist.* **39** 887–930. MR2816342 <https://doi.org/10.1214/10-AOS860>
- SUN, R. and LUO, Z.-Q. (2016). Guaranteed matrix completion via non-convex factorization. *IEEE Trans. Inf. Theory* **62** 6535–6579. MR3565131 <https://doi.org/10.1109/TIT.2016.2598574>
- SUN, T. and ZHANG, C.-H. (2012). Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *Adv. Neural Inf. Process. Syst.* 863–871.
- XIA, D. and YUAN, M. (2021). Statistical inferences of linear forms for noisy matrix completion. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **83** 58–77. MR4220984 <https://doi.org/10.1111/rssb.12400>
- ZHU, Z., WANG, T. and SAMWORTH, R. J. (2022). High-dimensional principal component analysis with heterogeneous missingness. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **84** 2000–2031. MR4515564