

LASSO UNDER MULTI-WAY CLUSTERING: ESTIMATION AND POST-SELECTION INFERENCE

HAROLD CHIANG AND YUYA SASAKI

ABSTRACT. This paper studies high-dimensional regression models with lasso when data is sampled under multi-way clustering. First, we establish convergence rates for the lasso and post-lasso estimators. Second, we propose a novel inference method based on a post-double-selection procedure and show its asymptotic validity. Our procedure can be easily implemented with existing statistical packages. Simulation results demonstrate that the proposed procedure works well in finite sample. We illustrate the proposed method with a couple of empirical applications to development and growth economics.

1. INTRODUCTION

This paper studies a method of estimation and post-selection inference for regression parameters in high-dimensional linear models by lasso under multi-way clustering. The objective is motivated by recently increasing demands from applied economic research. On one hand, economists often use multi-way cluster sampled data. Examples include, but are not limited to, network data, matched employer-employee data, matched student-teacher data, scanner data where observations are double-indexed by stores and products, market share data where observations are double-indexed by market and products, and growth/development data where observations are double-indexed by ethnicity and geographical units – see Section 7 for specific applications of the last example. On the other hand, researchers also often use machine learning methods of estimation and inference for high-dimensional models in today’s big data environments. There are a number of useful methods in the literature that deal with each of these two issues (multi-way clustering and

Date: First arXiv version: May 6, 2019.

Code files are available upon request from the authors.

Key words and phrases. cluster robust standard errors, high dimensions, lasso, machine learning, multi-way clustering, post-selection inference.

JEL Classification: C21, C55.

high dimensionality) separately, but the existing methods do not seem to provide a solution to dealing with both of these practically relevant issues simultaneously. In this light, we present lasso under multi-way clustering, and propose a post-selection inference method for regression parameters under this sampling assumption.

In the important branch of the literature following the seminal work by Belloni, Chen, Chernozhukov and Hansen (2012), post-selection inference with lasso has been widely studied under various settings by Belloni, Chernozhukov and Hansen (2014)¹, Javanmard and Montanari (2014), van de Geer, Bühlmann, Ritov (2014), Zhang and Zhang (2014), Belloni, Chernozhukov and Kato (2015), and many others. For empirical researchers, lasso has become a powerful machine learning tool under data-rich environments. Most of the papers in this literature assume i.i.d. or independent sampling. In many empirical applications, it is sometimes more plausible to assume multi-way cluster sampling (e.g., network data, matched employer-employee data, and matched student-teacher data). Building upon Belloni, Chen, Chernozhukov and Hansen (2012), Belloni and Chernozhukov (2013) and Belloni, Chernozhukov and Hansen (2014), this paper generalizes lasso and post-double-selection procedure by allowing for multi-way cluster sampling. To our best knowledge, the present paper is the first in the literature of high-dimensional models to consider lasso under multi-way cluster sampling.

The influential work by Cameron, Gelbach and Miller (2011) proposes multi-way cluster-robust inference methods for linear and nonlinear regression models – also see Cameron and Miller (2015, Section V) for a survey. Formal analysis of asymptotic properties and bootstrap validity under multi-way clustering is studied by Menzel (2017) using the Aldous-Hoover representation – see Kallenberg (2005, Chapter 7) for example. Under the assumptions of separable exchangeability, the method of Menzel (2017) covers both degenerate and non-degenerate cases. Using the same representation, while focusing on the non-degenerate cases, Davezies, D’Haultfoeuille and Guyonvarch (2018) develop empirical process theory under multi-way cluster sampling which applies to a large class of econometric models. Building upon the asymptotic framework of these two papers, MacKinnon, Nielsen and Webb (2019) propose several wild bootstrap procedures for linear regression models, and examine their finite-sample performances under several different cluster sampling scenarios. In this

¹See also Belloni, Chernozhukov and Hansen (2011).

paper, we take advantage of the innovations by these preceding papers to develop a multi-way cluster-robust inference method for high-dimensional models. To our best knowledge, the present paper is the first in this literature on multi-way clustering to consider high-dimensional models.

The rest of this paper is organized as follows. Section 2 introduces the model. Section 3 presents an overview of the proposed methodology. Section 4 discusses a formal asymptotic theory. Section 5 presents an extension of the baseline results to cases of heterogeneous cluster sizes. Section 6 presents simulation studies. Section 7 presents an empirical illustration with development and growth economics. Section 8 concludes. The appendix contains mathematical proofs and auxiliary lemmas.

2. THE MODEL

Consider the high-dimensional regression model

$$Y_{ij} = D_{ij}\alpha + X'_{ij}\beta + R^Y_{ij} + \varepsilon_{ij}, \quad \mathbb{E}[\varepsilon_{ij}|D_{ij}, X_{ij}] = 0, \quad (2.1)$$

where Y_{ij} is an observed outcome variable, $(D_{ij}, X'_{ij})'$ is an observed vector of regressors, and R^Y_{ij} is an approximation error for the unit of observation with the double index (i, j) . We set α as a scalar parameter of interest. The dimension p of the nuisance parameter vector $\beta \in \mathbb{R}^p$ is potentially increasing in the sample size. Following the literature on high-dimensional post-selection inference (e.g., Belloni, Chernozhukov and Hansen, 2014), we also consider the auxiliary projection

$$D_{ij} = X_{ij}\gamma + R^D_{ij} + v_{ij}, \quad \mathbb{E}[v_{ij}|X_{ij}] = 0, \quad (2.2)$$

where R^D is an approximation error. The dimension p of the nuisance parameter vector $\gamma \in \mathbb{R}^p$ is the same as that of β , and is potentially increasing in the sample size.

In the absence of two-way clustering, the system (2.1)–(2.2) would be the same as the model considered in Belloni, Chernozhukov and Hansen (2014). We first consider two-way clustering where each cell contains one observation. Section 5 presents an extension to the case of heterogeneous cluster sizes.

3. OVERVIEW OF THE METHOD

In this section, we present an overview of the proposed method, namely estimation and post-selection inference. Formal theoretical justifications are discussed in Section 4.

A researcher observes a sample $\{(Y_{ij}, D_{ij}, X'_{ij}) \mid i \in \{1, \dots, N\}, j \in \{1, \dots, M\}\}$ of size NM . The estimation procedure consists of two steps. First, define the lasso estimates for (2.1) and (2.2) by

$$(\hat{\alpha}, \hat{\beta}')' = \underset{\alpha, \beta}{\operatorname{argmin}} \sum_{i=1}^N \sum_{j=1}^M (Y_{ij} - \alpha D_{ij} - X'_{ij} \beta)^2 + \lambda_1 \|(\alpha, \beta)'\|_1, \quad (3.3)$$

$$\text{and } \hat{\gamma} = \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^N \sum_{j=1}^M (D_{ij} - X'_{ij} \gamma)^2 + \lambda_2 \|\gamma\|_1, \quad (3.4)$$

respectively, for some regularization parameters λ_1 and λ_2 , valid choices of which are discussed in the statement of Theorem 1 ahead. Denote the supports of the lasso estimates by $\hat{I}_1 = \operatorname{support}(\hat{\beta})$ and $\hat{I}_2 = \operatorname{support}(\hat{\gamma})$, and let $\hat{I} = \hat{I}_1 \cup \hat{I}_2$. In the second step, define the post-double-selection lasso estimate $\tilde{\alpha}$ by

$$(\tilde{\alpha}, \tilde{\beta}') = \underset{\operatorname{support}(\beta) \subset \hat{I}}{\operatorname{argmin}} \sum_{i=1}^N \sum_{j=1}^M (Y_{ij} - \alpha D_{ij} - X'_{ij} \beta)^2. \quad (3.5)$$

Let $\underline{C} = N \wedge M$, $\mu_N = \underline{C}/N$ and $\mu_M = \underline{C}/M$. Under suitable conditions to be formally stated in Section 4, we have the asymptotic normality

$$\sigma^{-1} \sqrt{\underline{C}} (\tilde{\alpha} - \alpha) \rightsquigarrow N(0, 1),$$

where the asymptotic variance is given by $\sigma^2 = Q^{-1} \Gamma Q^{-1}$ with

$$Q = \mathbb{E}[v_{11}^2],$$

$$\Gamma = \bar{\mu}_N \Gamma_N + \bar{\mu}_M \Gamma_M = \bar{\mu}_N \mathbb{E}[v_{11} \varepsilon_{11} v_{12} \varepsilon_{12}] + \bar{\mu}_M \mathbb{E}[v_{11} \varepsilon_{11} v_{21} \varepsilon_{21}],$$

and $\bar{\mu}_N$ and $\bar{\mu}_M$ denoting the limits of μ_N and μ_M , respectively.

The asymptotic variance is estimated by the sample counterpart $\hat{\sigma}^2 = \hat{Q}^{-1}\hat{\Gamma}\hat{Q}^{-1}$, where

$$\hat{Q} = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M \hat{v}_{ij}^2,$$

$$\hat{\Gamma} = \frac{\underline{C}}{(NM)^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} \hat{v}_{ij} \hat{\varepsilon}_{ij} \hat{\varepsilon}_{ij'} \hat{v}_{ij'} + \frac{\underline{C}}{(NM)^2} \sum_{1 \leq i, i' \leq N} \sum_{j=1}^M \hat{v}_{ij} \hat{\varepsilon}_{ij} \hat{\varepsilon}_{i'j} \hat{v}_{i'j},$$

$$\hat{v}_{ij} = D_{ij} - X'_{ij}\hat{\gamma}, \text{ and } \hat{\varepsilon}_{ij} = Y_{ij} - \hat{\alpha}D_{ij} - X'_{ij}\hat{\beta}.$$

In summary, we propose to report the post-double-selection lasso estimate $\tilde{\alpha}$ as an estimate of α with its standard error given by $\hat{\sigma}/\sqrt{\underline{C}}$. The α^* -level confidence interval can be constructed as $[\tilde{\alpha} + \Phi^{-1}(\alpha^*/2)\hat{\sigma}/\sqrt{\underline{C}}, \tilde{\alpha} + \Phi^{-1}(1 - \alpha^*/2)\hat{\sigma}/\sqrt{\underline{C}}]$, where Φ^{-1} denotes the quantile function of the standard normal distribution.

4. ASYMPTOTIC THEORY

The two-way sample sizes $(N, M) \in \mathbb{N}^2$ will be indexed by a single index $n \in \mathbb{N}$ as $(N, M) = (N(n), M(n))$ where $M(n)$ and $N(n)$ are non-decreasing in n and $M(n)N(n)$ is increasing in n . For simplicity, each size of intersection n_{ij} is assumed to be uniformly bounded by a positive integer \bar{n} that is independent of n . With this said, we will suppress the index notation and write (N, M) for simplicity. We fix a number of notations. For each n , let \mathbb{P}_n denote the law with respect to sample size (N, M) – note that we allow the dimension p of X_{ij} to grow with n . Let $a := p \vee (NM)$. Also recall the notations $\underline{C} = N \wedge M$, $\mu_N = \underline{C}/N$, and $\mu_M = \underline{C}/M$ from Section 3. We use the short-hand notation $[k] = \{1, \dots, k\}$ and $[k]^c = \mathbb{N} \setminus [k]$ for any $k \in \mathbb{N}$. For a sequence $(t_{ij})_{i \in [N], j \in [M]}$, denote $\|t_{ij}\|_n = \sqrt{\frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M t_{ij}^2}$. Thus, $\|X'_{ij}\delta\|_n = \sqrt{\frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M \delta' X_{ij} X'_{ij} \delta}$ is the prediction norm of δ . Let $\|A\|_\infty = \max_{k,l} |A_{k,l}|$ denote the max norm of matrix A . We write $a \lesssim b$ to mean $a \leq cb$ for some $c > 0$ that does not depend on n . We also write $a \lesssim_P b$ to mean $a = O_P(b)$. We write $Z_{ij} = (Y_{ij}, D_{ij}, X'_{ij})'$ for the $(p+2)$ -dimensional random vector in data. Throughout, we assume that this random vector Z_{ij} is Borel measurable – see Kallenberg (2005, pp. 304). With these notations, we state the following four assumptions.

Assumption 1 (Sampling). *Suppose that $\underline{C} \rightarrow \infty$, $\mu_N \rightarrow \bar{\mu}_N \geq 0$, and $\mu_M \rightarrow \bar{\mu}_M \geq 0$.*

- (1) $(Z_{ij})_{(i,j) \in \mathbb{N}^2}$ is an infinite sequence of separately exchangeable $(p+2)$ -dimensional random vectors. That is, for any permutations π_1 and π_2 of \mathbb{N} , we have

$$(Z_{ij})_{(i,j) \in \mathbb{N}^2} \stackrel{d}{=} (Z_{\pi_1(i)\pi_2(j)})_{(i,j) \in \mathbb{N}^2}.$$

- (2) $(Z_{ij})_{(i,j) \in \mathbb{N}^2}$ is dissociated. That is, for any $(c_1, c_2) \in \mathbb{N}^2$, $(Z_{ij})_{i \in [c_1], j \in [c_2]}$ is independent of $(Z_{ij})_{i \in [c_1]^c, j \in [c_2]^c}$.
- (3) For each n , an econometrician observes $(Z_{ij})_{i \in [N], j \in [M]}$.

Assumption 2 (Moments). *There exists a sequence $\{B_n\}_{n=1}^\infty$ of positive constants such that the following conditions hold for all $n \in \mathbb{N}$ for some $q > 4$:*

- (1) $\mathbb{E}[|D_{11}|^{2q}] + \max_{k \in [p]} \mathbb{E}[|X_{11,k}|^{2q}] + \mathbb{E}[|\varepsilon_{11}|^{2q}|X_{11}, v_{11}] + \mathbb{E}[|v_{11}|^{2q}|X_{11}] \leq K$ a.s. and $0 < c \leq \mathbb{E}[v_{11}^2|X_{11}]$ a.s. for positive constants, c and K , that are independent of n .
- (2) $\mathbb{E}[\|X_{11}\|_\infty^{2q}] \leq B_n^{2q}$ and $B_n \sqrt{\log a} \lesssim (N \vee M)^{1/2-1/q}$.
- (3) $\bar{\mu}_N \mathbb{E}[v_{11}\varepsilon_{11}v_{12}\varepsilon_{12}] + \bar{\mu}_M \mathbb{E}[v_{11}\varepsilon_{11}v_{21}\varepsilon_{21}]$ and the maximal and minimal eigenvalues of $\mathbb{E}[X_{11}X'_{11}]$ are bounded and bounded away from zero uniformly in n .

Assumption 3 (Sparsity).

- (1) $\|\beta\|_0 + \|\gamma\|_0 \lesssim s$.
- (2) $\|R_{ij}^Y\|_n + \|R_{ij}^D\|_n \leq c_s \lesssim_P \sqrt{s/\underline{C}}$.
- (3) $\frac{s_n^2(\log(a))^2}{\underline{C}} = o(1)$.

Assumption 4 (Sparse Eigenvalues). *There exists a sequence $\{\ell_n\}$ such that $\ell_n \rightarrow \infty$ and, with probability at least $1 - o(1)$,*

$$0 < c \leq \phi_{\min}(s\ell_n) \leq \phi_{\max}(s\ell_n) \leq c' < \infty$$

holds for some constants, c and c' , that are independent of n , where

$$\phi_{\max}(m) := \max_{1 \leq \|\delta\|_0 \leq m} \frac{\delta' M \delta}{\|\delta\|^2} \text{ and } \phi_{\min}(m) := \min_{1 \leq \|\delta\|_0 \leq m} \frac{\delta' M \delta}{\|\delta\|^2},$$

with $M = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M X_{ij}X'_{ij}$, denote the maximal and minimal m -sparse eigenvalues.

Remark 4.1 (Discussion of the Assumptions). Assumption 1 is closely related to Assumption 1 of Davezies, D'Haultfoeuille and Guyonvarch (2018). The main difference is that we allow p to be changing with n . We remark that the exchangeability assumption is not new in econometrics – it has been used in Andrews (2005) and Menzel (2015) as well as Menzel

(2017), Davezies, D'Haultfoeuille and Guyonvarch (2018), and MacKinnon, Nielsen and Webb (2019). Assumption 2 is standard in the literature on post-selection inference with lasso. Parts (1) and (2) require an existence of higher order moments of key objects. Note that common assumptions in high-dimensional literature, such as sub-gaussianity or boundedness, are not required. They can be replaced by some higher level conditions similar to Condition RF of Belloni, Chen, Chernozhukov and Hansen (2012).² Part (3) of Assumption 2 requires that the asymptotic variance is bounded away from zero.³ Assumption 3 is a direct generalization of Condition ASTE (iii) and (iv) of Belloni, Chernozhukov and Hansen (2014). Finally, Assumption 4 is analogous to Condition SE of Belloni, Chernozhukov and Hansen (2014), which is standard in the high-dimensional literature. It only imposes small diagonal submatrices to be well behaved. \triangle

4.1. Independentization via Hájek Projection. In this section, we show that an empirical process in multi-way clustered samples can be represented as a sum of independent variables via Hájek projection. Furthermore, its variance can be shown to be approximated by covariances of observed variables.

For any $f : \text{support}(Z) \rightarrow \mathbb{R}$, we let

$$\mathbb{G}_C f := \sqrt{C} \left\{ \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M f(Z_{ij}) - \mathbb{E}[f(Z_{11})] \right\}$$

denote its empirical process.

Lemma 1 (Independentization via Hájek Projection). *If Assumption 1 holds and $f : \text{support}(Z) \rightarrow \mathbb{R}$ satisfies $\mathbb{E}f^2(Z_{11}) < K$ for a finite constant K that is independent of n , then there exist i.i.d. uniform random variables U_{i0} and U_{0j} such that the Hájek projection $H_n f$ of $\mathbb{G}_C f$ on*

$$\mathcal{G}_n = \left\{ \sum_{i=1}^N g_{i0}(U_{i0}) + \sum_{j=1}^M g_{0j}(U_{0j}) : g_{i0}, g_{0j} \in L^2(\mathbb{P}_n) \right\}$$

²See their Lemma 3.

³Similarly to Davezies, D'Haultfoeuille and Guyonvarch (2018), we focus on non-degenerate cases in this paper. See Menzel (2017) for the studies of degenerate cases using a bootstrap-based method.

is equal to

$$H_n f = \sum_{i=1}^N \frac{\sqrt{\underline{C}}}{N} \mathbb{E} \left[f(Z_{i1}) - \mathbb{E} f(Z_{11}) \middle| U_{i0} \right] + \sum_{j=1}^M \frac{\sqrt{\underline{C}}}{M} \mathbb{E} \left[f(Z_{1j}) - \mathbb{E} f(Z_{11}) \middle| U_{0j} \right]$$

for each n . Furthermore,

$$V(\mathbb{G}_C f) = V(H_n f) + O(\underline{C}^{-1}) = \bar{\mu}_N \text{Cov}(f(Z_{11}), f(Z_{12})) + \bar{\mu}_M \text{Cov}(f(Z_{11}), f(Z_{21})) + O(\underline{C}^{-1})$$

holds a.s.

A proof of this lemma can be found in Appendix A.1. The first part of the lemma shows that an empirical process $\mathbb{G}_C f$ under multi-way cluster sampling can be represented as a sum of independent unobserved variables via Hájek projection $H_n f$. While U_{i0} and U_{0j} are unobserved, the second part of this lemma in turn shows that the variance of the Hájek projection can be approximated by covariances of observed variables. Note that, since $H_n f$ is a Hájek projection, the lemma implies $\frac{\mathbb{G}_C f}{\sqrt{V(\mathbb{G}_C f)}} = \frac{H_n f}{\sqrt{V(H_n f)}} + o_P(1)$ if $\bar{\mu}_N \text{Cov}(f(Z_{11}), f(Z_{12})) + \bar{\mu}_M \text{Cov}(f(Z_{11}), f(Z_{21}))$ is bounded and bounded away from zero uniformly in n .

Our Lemma 1 can be seen as an extension to Lemma D.2 in Davezies, D'Haultfoeuille and Guyonvarch (2018). Specifically, while Davezies, D'Haultfoeuille and Guyonvarch (2018) consider a fixed data generating process over the sample size n , our Lemma 1 allows the data generating process to vary with n in particular for the sake of accommodating the increasing of dimensionality p for high-dimensional models. The lemma serves as a main building block for all the asymptotic results to be presented ahead.

4.2. Convergence Rates of Lasso and Post-Lasso under Multi-Way Clustering.

We next show the convergence rates of the lasso estimator $(\hat{\alpha}, \hat{\beta}', \hat{\gamma}')'$ and the post-lasso estimator $(\tilde{\alpha}, \tilde{\beta}', \tilde{\gamma}')'$ under multi-way clustering.

Theorem 1 (Convergence Rates for Lasso and Post-Lasso under Multi-Way Clustering). *If Assumptions 1, 2 (1)–(2), 3 (1)–(2), and 4 are satisfied, and $\lambda_1, \lambda_2 = C \sqrt{(NM)^2 \log a / \underline{C}}$*

for some constant $C > 1$, then

$$\begin{aligned} \|\hat{\eta} - \eta\|_1 + \|\hat{\gamma} - \gamma\|_1 &\lesssim \sqrt{\frac{s^2 \log a}{\underline{C}}}, & \|W'_{ij}(\hat{\eta} - \eta)\|_n + \|X_{ij}(\hat{\gamma} - \gamma)\|_n &\lesssim \sqrt{\frac{s \log a}{\underline{C}}}, \\ \|\tilde{\eta} - \eta\|_1 + \|\tilde{\gamma} - \gamma\|_1 &\lesssim \sqrt{\frac{s^2 \log a}{\underline{C}}}, & \|W'_{ij}(\tilde{\eta} - \eta)\|_n + \|X_{ij}(\tilde{\gamma} - \gamma)\|_n &\lesssim \sqrt{\frac{s \log a}{\underline{C}}}, \quad \text{and} \\ \|\hat{\eta} - \eta\| + \|\hat{\gamma} - \gamma\| + \|\tilde{\eta} - \eta\| + \|\tilde{\gamma} - \gamma\| &\lesssim \sqrt{\frac{s \log a}{\underline{C}}} \end{aligned}$$

hold, where $W_{ij} = [D_{ij}, X'_{ij}]'$ and $\eta = (\alpha, \beta)'$.

A proof can be found in Appendix A.2, and is based on the previous result (Lemma 1). In the multi-way sampling, this lemma can be viewed as a counterpart of Lemma 6 and Lemma 7 in Belloni, Chen, Chernozhukov and Hansen (2012).

4.3. Post-Selection-Inference with Post-Lasso under Multi-way Clustering. In this section, we present the main result of this paper. The limit normal distribution of the post-double-selection lasso estimate $\tilde{\alpha}$ is established based on the previous two results (Lemma 1 and Theorem 1).

Theorem 2 (Asymptotic Normality). *If Assumptions 1, 2, 3 and 4 are satisfied, and λ_1 and λ_2 are chosen according to the statement of Theorem 1, then*

$$\sigma^{-1} \sqrt{\underline{C}}(\tilde{\alpha} - \alpha) \rightsquigarrow N(0, 1),$$

where $\sigma^2 = Q^{-1} \Gamma Q^{-1}$, $Q = \mathbb{E}[v_{11}^2]$ and

$$\Gamma = \bar{\mu}_N \Gamma_N + \bar{\mu}_M \Gamma_M = \bar{\mu}_N \mathbb{E}[v_{11} \varepsilon_{11} v_{12} \varepsilon_{12}] + \bar{\mu}_M \mathbb{E}[v_{11} \varepsilon_{11} v_{21} \varepsilon_{21}].$$

A proof can be found in Appendix A.3. This result provides a theoretical justification for the asymptotic variance proposed in the overview in Section 3. In practice, we do not know the components, Q and Γ , of the asymptotic variance. The following subsection proposes estimators of them.

4.4. Variance Estimation. In this section, we propose an analog variance estimator. The components, Q and Γ , of the asymptotic variance can be estimated by

$$\widehat{Q} = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M \widehat{v}_{ij}^2 \quad \text{and}$$

$$\widehat{\Gamma} = \frac{\underline{C}}{(NM)^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} \widehat{v}_{ij} \widehat{\varepsilon}_{ij} \widehat{\varepsilon}_{ij'} \widehat{v}_{ij'} + \frac{\underline{C}}{(NM)^2} \sum_{1 \leq i, i' \leq N} \sum_{j=1}^M \widehat{v}_{ij} \widehat{\varepsilon}_{ij} \widehat{\varepsilon}_{i'j} \widehat{v}_{i'j},$$

respectively, where $\widehat{v}_{ij} = D_{ij} - X'_{ij} \widehat{\gamma}$ and $\widehat{\varepsilon}_{ij} = Y_{ij} - \widehat{\alpha} D_{ij} - X'_{ij} \widehat{\beta}$ are the residuals. With these component estimators, we propose that the asymptotic variance $\sigma^2 = Q^{-1} \Gamma Q^{-1}$ be estimated by $\widehat{\sigma}^2 = \widehat{Q}^{-1} \widehat{\Gamma} \widehat{Q}^{-1}$. The following theorem provides a theoretical support for this variance estimator.

Theorem 3 (Variance Estimation). *If Assumptions 1, 2, 3 and 4 are satisfied, λ_1 and λ_2 are chosen according to the statement of Theorem 1, $\frac{(NM)^{1/q} B_n^2 s^3 (\log a)^2}{\underline{C}^2} = o(1)$, $\frac{(NM)^{1/q} s \log a}{\underline{C}} = o(1)$, and $\|R_{ij}^D R_{ij}^Y\|_n^2 = O(1)$, then the variance estimator $\widehat{\sigma}^2 = \widehat{Q}^{-1} \widehat{\Gamma} \widehat{Q}^{-1}$ is consistent for $\sigma^2 = Q^{-1} \Gamma Q^{-1}$.*

A proof is found in Appendix A.4. In light of this result, we propose to compute the standard error by $\widehat{\sigma}/\sqrt{\underline{C}}$. Similarly, in light of this result together with Theorem 2, we propose to construct the α^* -level confidence interval by $[\widetilde{\alpha} + \Phi^{-1}(\alpha^*/2) \widehat{\sigma}/\sqrt{\underline{C}}, \widetilde{\alpha} + \Phi^{-1}(1 - \alpha^*/2) \widehat{\sigma}/\sqrt{\underline{C}}]$, where Φ^{-1} denotes the quantile function of the standard normal distribution.

5. EXTENSION: HETEROGENEOUS CLUSTER SIZES

Thus far, we focus on the case where each cluster contains one observation. In this section, we presented an extension of the baseline results to situations where the numbers of observations are heterogeneous across clusters. Suppose that we have n_{ij} observations for each cell $(i, j) \in [N] \times [M]$, where n_{ij} is a random variable that is allowed to depend on $(X_{ij,\ell})_{\ell \geq 1}$. To deal with the situation of $n_{ij} = 0$, for any sequence $(t_\ell)_{\ell \geq 1}$, define $\sum_{\ell=1}^0 t_\ell = 0$. Consider the model

$$Y_{ij,\ell} = D_{ij,\ell} \alpha + X'_{ij,\ell} \beta + R_{ij,\ell}^Y + \varepsilon_{ij,\ell}, \quad \mathbb{E}[\varepsilon_{ij,\ell} | D_{ij,\ell}, X_{ij,\ell}] = 0,$$

where $Y_{ij,\ell}$ is an observed outcome variable, $(D_{ij,\ell}, X'_{ij,\ell})'$ is an observed vector of regressors, and $R_{ij,\ell}^Y$ is an approximation error for the unit $\ell \in [n_{ij}]$ with the double index (i, j) indicating i -th cluster in the first clustering dimension and j -th cluster in the second clustering dimension. Using matrix notations, we can rewrite the model as

$$Y_{ij} = D_{ij}\alpha + X_{ij}\beta + R_{ij}^Y + \varepsilon_{ij}, \quad \text{E}[\varepsilon_{ij}|D_{ij}, X_{ij}] = 0,$$

where each of $Y_{ij} = (Y_{ij,\ell})_{\ell \in [n_{ij}]}$, $D_{ij} = (D_{ij,\ell})_{\ell \in [n_{ij}]}$, $R_{ij}^Y = (R_{ij,\ell}^Y)_{\ell \in [n_{ij}]}$, and $\varepsilon_{ij} = (\varepsilon_{ij,\ell})_{\ell \in [n_{ij}]}$ is of dimension $n_{ij} \times 1$, and $X_{ij} = (X'_{ij,\ell})_{\ell \in [n_{ij}]}$ is of dimension $n_{ij} \times p$. We similarly write the accompanying auxiliary projection as

$$D_{ij} = X_{ij}\gamma + R_{ij}^D + v_{ij}, \quad \text{E}[v_{ij}|X_{ij}] = 0,$$

where R^D is of dimension $n_{ij} \times 1$ representing approximation errors, and v_{ij} is of dimension $n_{ij} \times 1$ representing projection errors.

Under this setting, the first step of estimation procedure consists of

$$\begin{aligned} (\hat{\alpha}, \hat{\beta})' &= \underset{\alpha, \beta}{\operatorname{argmin}} \sum_{i=1}^N \sum_{j=1}^M \sum_{\ell=1}^{n_{ij}} (Y_{ij,\ell} - \alpha D_{ij,\ell} - X'_{ij,\ell}\beta)^2 + \lambda_1 \|(\alpha, \beta)'\|_1 \\ \text{and } \hat{\gamma} &= \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^N \sum_{j=1}^M \sum_{\ell=1}^{n_{ij}} (D_{ij,\ell} - X'_{ij,\ell}\gamma)^2 + \lambda_2 \|\gamma\|_1. \end{aligned}$$

In turn, the second-step estimates are obtained by

$$(\tilde{\alpha}, \tilde{\beta}') = \underset{\operatorname{support}(\beta) \subset \hat{I}}{\operatorname{argmin}} \sum_{i=1}^N \sum_{j=1}^M \sum_{\ell=1}^{n_{ij}} (Y_{ij,\ell} - \alpha D_{ij,\ell} - X'_{ij,\ell}\beta)^2.$$

The asymptotic variance estimator for $\tilde{\alpha}$ is given by $\hat{\sigma}^2 = \hat{Q}^{-1}\hat{\Gamma}\hat{Q}^{-1}$, where

$$\begin{aligned} \hat{Q} &= \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M \hat{v}'_{ij} \hat{v}_{ij}, \\ \hat{\Gamma} &= \frac{\underline{C}}{(NM)^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} \hat{v}'_{ij} \hat{\varepsilon}_{ij} \hat{\varepsilon}'_{ij'} \hat{v}_{ij'} + \frac{\underline{C}}{(NM)^2} \sum_{1 \leq i, i' \leq N} \sum_{j=1}^M \hat{v}'_{ij} \hat{\varepsilon}_{ij} \hat{\varepsilon}'_{i'j} \hat{v}_{i'j}, \end{aligned}$$

$$\hat{v}_{ij} = D_{ij} - X_{ij}\hat{\gamma}, \quad \text{and } \hat{\varepsilon}_{ij} = Y_{ij} - \hat{\alpha}D_{ij} - X_{ij}\hat{\beta}.$$

We now formally state assumptions for the extended theory to support the asymptotic validity of this procedure. Define $W_{ij} = (n_{ij}, (Z_{ij,\ell})_{\ell \geq 1})$ and $\ddot{M} = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M X'_{ij} X_{ij}$.

Assumption 5 (Sampling). *Suppose that $\underline{C} \rightarrow \infty$, $\mu_N \rightarrow \bar{\mu}_N \geq 0$, and $\mu_M \rightarrow \bar{\mu}_M \geq 0$.*

- (1) $(W_{ij})_{(i,j) \in \mathbb{N}^2}$ is an infinite sequence of separately exchangeable random processes.
- (2) $(W_{ij})_{(i,j) \in \mathbb{N}^2}$ is dissociated.
- (3) For each n , an econometrician observes $((W_{ij,\ell})_{\ell \in [n_{ij}]})_{i \in [N], j \in [M]}$.
- (4) $\mathbb{E}[n_{ij}] > 0$ and $n_{ij} \leq \bar{n}$ for a positive finite constant \bar{n} independent of n .

Assumption 6 (Moments). *There exists a sequence $\{B_n\}_{n=1}^\infty$ of positive constants such that the following conditions hold for all $n \in \mathbb{N}$ for some $q > 4$:*

- (1) $\mathbb{E}[\max_{\ell \in [n_{ij}]} |D_{11,\ell}|^{2q}] + \max_{k \in [p]} \mathbb{E}[\max_{\ell \in [n_{ij}]} |X_{11,\ell,k}|^{2q}] + \mathbb{E}[\max_{\ell \in [n_{ij}]} |\varepsilon_{11,\ell}|^{2q} | X_{11,\ell}, v_{11,\ell}] + \mathbb{E}[|v_{11,\ell}|^{2q} | X_{11,\ell}] \leq K$ a.s. and $0 < c \leq \mathbb{E}[\max_{\ell \in [n_{ij}]} v_{11,\ell}^2 | X_{11,\ell}]$ a.s. for positive constants, c and K , that are independent of n .
- (2) $\mathbb{E}[\max_{\ell \in [n_{ij}]} \|X_{11}\|_\infty^{2q}] \leq B_n^{2q}$ and $B_n \sqrt{\log a} \lesssim (N \vee M)^{1/2-1/q}$.
- (3) $\bar{\mu}_N \mathbb{E}[v'_{11} \varepsilon_{11} \varepsilon'_{12} v_{12}] + \bar{\mu}_M \mathbb{E}[v'_{11} \varepsilon_{11} \varepsilon'_{21} v_{21}]$ and the maximal and minimal eigenvalues of $\mathbb{E}[X'_{11} X_{11}]$ are bounded and bounded away from zero uniformly in n .

Assumption 7 (Sparsity).

- (1) $\|\beta\|_0 + \|\gamma\|_0 \lesssim s$.
- (2) $\sqrt{(NM)^{-1} \sum_{i=1}^N \sum_{j=1}^M \sum_{\ell \in [n_{ij}]} (R_{ij,\ell}^Y)^2} + \sqrt{(NM)^{-1} \sum_{i=1}^N \sum_{j=1}^M \sum_{\ell \in [n_{ij}]} (R_{ij,\ell}^D)^2} \leq c_s \lesssim_P \sqrt{s/\underline{C}}$.
- (3) $\frac{s_n^2 (\log(a))^2}{\underline{C}} = o(1)$.

Assumption 8 (Sparse Eigenvalues). *There exists a sequence $\{\ell_n\}$ such that $\ell_n \rightarrow \infty$ and, with probability at least $1 - o(1)$,*

$$0 < c \leq \phi_{\min}(s\ell_n) \leq \phi_{\max}(s\ell_n) \leq c' < \infty$$

holds for some constants, c and c' , that are independent of n , where

$$\phi_{\max}(m) := \max_{1 \leq \|\delta\|_0 \leq m} \frac{\delta' \ddot{M} \delta}{\|\delta\|^2} \text{ and } \phi_{\min}(m) := \min_{1 \leq \|\delta\|_0 \leq m} \frac{\delta' \ddot{M} \delta}{\|\delta\|^2}.$$

The following statement provides a theoretical guarantee for the estimation and inference procedure for the extended model outlined above.

Corollary 1. *If Assumptions 5, 6, 7 and 8 are satisfied, and λ_1 and λ_2 are chosen according to the statement of Theorem 1, then*

$$\sigma^{-1} \sqrt{\underline{C}} (\tilde{\alpha} - \alpha) \rightsquigarrow N(0, 1),$$

where $\sigma^2 = Q^{-1}\Gamma Q^{-1}$ with

$$Q = \mathbb{E}[v'_{11}v_{11}],$$

$$\Gamma = \bar{\mu}_N \mathbb{E}[v'_{11}\varepsilon_{11}v'_{12}\varepsilon_{12}] + \bar{\mu}_M \mathbb{E}[v'_{11}\varepsilon_{11}v'_{21}\varepsilon_{21}].$$

Furthermore, if $\frac{(NM)^{1/q}B_n^2s^3(\log a)^2}{\underline{C}^2} = o(1)$, $\frac{(NM)^{1/q}s \log a}{\underline{C}} = o(1)$, and $\|R_{ij}^D R_{ij}^Y\|_n^2 = O(1)$, then the variance estimator $\hat{\sigma}^2$ is consistent for σ^2 .

A proof of Corollary 1 closely follows that of the results in Section 4, and are therefore omitted. The key difference is that we now apply Aldous-Hoover representation on W_{ij} rather than on Z_{ij} .⁴

6. SIMULATION STUDIES

In this section, we present simulation studies of finite-sample performance of the proposed method of estimation and post-selection inference. We compare the performance of our method against existing alternatives from the lasso literature that do not account for multi-way clustering.

6.1. Simulation Setup. We consider the linear model

$$Y_{ij} = D_{ij}\alpha + X'_{ij}\beta + \varepsilon_{ij}.$$

The parameter values are fixed at $(\alpha, \beta)' = (0.5, 0.5^2, \dots, 0.5^{\dim(X)+1})'$. The random vector $(D_{ij}, X'_{ij}, \varepsilon_{ij})$ is constructed by

$$(D_{ij}, X_{ij}) = (1 - \omega_1^x - \omega_2^x)v_{ij}^x + \omega_1^x v_i^x + \omega_2^x v_j^x \quad \text{and}$$

$$\varepsilon_{ij} = (1 - \omega_1^\varepsilon - \omega_2^\varepsilon)v_{ij}^\varepsilon + \omega_1^\varepsilon v_i^\varepsilon + \omega_2^\varepsilon v_j^\varepsilon$$

⁴For more insights on this extension, see Section 3.1 of Davezies, D'Haultfoeuille and Guyonvarch (2019).

with two-way clustering weights (ω_1^x, ω_2^x) and $(\omega_1^\varepsilon, \omega_2^\varepsilon)$, where v_{ij}^x , v_i^x , and v_j^x are independently generated according to

$$v_{ij}^x, v_i^x, v_j^x \sim N \left(0, \begin{pmatrix} \rho^0 & \rho^1 & \dots & \rho^{\dim(X)-1} & \rho^{\dim(X)} \\ \rho^1 & \rho^0 & \dots & \rho^{\dim(X)-2} & \rho^{\dim(X)-1} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \rho^{\dim(X)-1} & \rho^{\dim(X)-2} & \dots & \rho^0 & \rho^1 \\ \rho^{\dim(X)} & \rho^{\dim(X)-1} & \dots & \rho^1 & \rho^0 \end{pmatrix} \right),$$

and v_{ij}^ε , v_i^ε , and v_j^ε are independently generated according to

$$v_{ij}^\varepsilon, v_i^\varepsilon, v_j^\varepsilon \sim N(0, 1).$$

Note that the weights (ω_1^x, ω_2^x) and $(\omega_1^\varepsilon, \omega_2^\varepsilon)$ specify the extent of dependence in two-way clustering in (D_{ij}, X'_{ij}) and ε_{ij} , respectively. Also, the parameter ρ specifies the extent of collinearity among the high-dimensional covariates (D_{ij}, X'_{ij}) . We set $(\omega_1^x, \omega_2^x) = (0.25, 0.25)$, $(\omega_1^\varepsilon, \omega_2^\varepsilon) = (0.25, 0.25)$, and $\rho = 0.50$.

6.2. Alternative Variance Estimators. We compare the performance of our multi-way cluster-robust variance estimator with two existing alternative benchmarks. One is the heteroskedasticity robust variance estimator (such as the one in Belloni, Chernozhukov and Hansen (2014)) without accounting for cluster sampling, i.e., Γ is estimated by

$$\widehat{\Gamma}_{HC} = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M \widehat{v}_{ij}^2 \widehat{\varepsilon}_{ij}^2.$$

We will refer to this variance estimator $\widehat{Q}^{-1} \widehat{\Gamma}_{HC} \widehat{Q}^{-1}$ as the ‘0-Way’ estimator. The other is the one-way cluster-robust variance estimator (similar to those of Belloni, Chernozhukov and Hansen and Kozlowski (2016) and Kock (2016)) clustered at one (e.g., second) dimension, i.e., Γ is estimated by

$$\widehat{\Gamma}_{CR} = \frac{1}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} \widehat{v}_{ij} \widehat{\varepsilon}_{ij} \widehat{v}_{ij'} \widehat{\varepsilon}_{ij'}.$$

We will refer to this variance estimator $\widehat{Q}^{-1} \widehat{\Gamma}_{CR} \widehat{Q}^{-1}$ as the ‘1-Way’ estimator.

6.3. Results. Table 1 summarizes simulation results. The first two columns indicate the two-way sample sizes (N, M) . The third column indicates the dimension (Dim) of $(\alpha, \beta)'$. The next four columns report simulation statistics for $\tilde{\alpha}$. These statistics include the average (Avg), bias (Bias), standard deviation (SD), and root mean square error (RMSE). The last three columns report 95% coverage frequencies of α based on three variance estimators. The first is the heteroskedasticity robust variance estimator (0-Way). The second is the one-way cluster-robust variance estimator (1-Way). The third is our multi-way cluster-robust variance estimator (2-Way). The results are based on 25,000 Monte Carlo iterations for each row in the table.

In view of the statistics columns, observe that the post-double-selection lasso estimate $\tilde{\alpha}$ behaves well in larger sample sizes (e.g., $N, M \geq 20$) both in terms of bias and variance. Next, observe the 95% coverage frequencies by the three alternative variance estimators. Both the 0-Way and 1-Way variance estimators significantly underestimate the variances of the post-double-selection lasso estimate $\tilde{\alpha}$. On the other hand, the coverage frequency based on our 2-Way variance estimator approaches the nominal probability (95%) as the sample size increases. These results demonstrate that, when the true sampling process entails multi-way clustering, traditional variance estimators may bias the inference and our multi-way cluster-robust variance estimator performs robustly well.

7. EMPIRICAL ILLUSTRATIONS

In this section, we illustrate our proposed method with applications to a couple of empirical studies. There is a sequence of recent growth and development economic studies using empirical data that are clustered at ethnic and geographical levels (e.g., Nunn and Wantchekon, 2011; Michalopoulos and Papaioannou, 2013, 2014, 2016; Gershman, 2016; Anderson, 2018; Dickens, 2018). The next two subsections present how our method can enrich the model flexibility and robustness of such studies, focusing on the cases of Nunn and Wantchekon (2011) and Michalopoulos and Papaioannou (2013).

7.1. Slave Trade and Mistrust in Africa. Nunn and Wantchekon (2011) analyze the effects of slave trade on mistrust in Africa, controlling for various demographic and geographical covariates including age, age squared, ethnic fractionalization, gender, urban residence, occupation, religion, and living conditions as well as country fixed effects in their

baseline model. Estimates of these effects are obtained by running regressions with a sample that pools n_{ij} individuals $\ell \in [n_{ij}]$ in ethnic group i and districts j across the cells $(i, j) \in [N] \times [M]$ of $N(= 185)$ ethnic groups and $M(= 1257)$ districts. Standard errors are computed by the two-way cluster-robust method of Cameron, Gelbach and Miller (2011) for the ethnic group and district as two ways of clustering.

With our proposed method that is applicable to both high-dimensional models and multi-way clustering, they could consider even more flexible model specifications, for example, allowing for higher orders of age rather than just the quadratic specification and interactions of the age polynomials with various other dummy variables. We present estimates with standard errors under such extended models with flexible specifications, demonstrate that qualitatively similar results continue to be obtained without substantial loss of statistical significance, and thus confirm further robustness of the main empirical findings by Nunn and Wantchekon (2011).

Consider the model

$$Y_{ij,\ell} = D_{ij,\ell}\alpha + X'_{ij,\ell}\beta + R_{ij,\ell}^Y + \varepsilon_{ij,\ell}, \quad \text{E}[\varepsilon_{ij,\ell}|D_{ij,\ell}, X_{ij,\ell}] = 0,$$

where $Y_{ij,\ell}$ denotes a measure of trust, $D_{ij,\ell}$ denotes an intensity measure of slave trade, $X_{ij,\ell}$ contains polynomial basis elements of age up to degree 10, ethnic fractionalization, gender, urban residence, occupation, religion, living conditions, the interactions of the polynomial basis of age with all the dummy variables, and country fixed effects, consisting of 597 dimensions of covariates in total. Note that the total number of regressors ($p + 1 = 598$) is much larger than the effective sample size ($\underline{C} = N \wedge M = 185$) of two-way clustering in this extended setting.

Table 2 summarize the estimates of the effects of slave trade on mistrust as measured by the “trust of neighbors,” corresponding to Table 1 of Nunn and Wantchekon (2011). The last two columns in the table show the original estimates obtained under the prototypical model by Nunn and Wantchekon (2011, Table 1) and corresponding lasso estimates obtained under more flexible model specification by our method. Across all the measures of slave exports, the original estimates and our lasso estimates are similar with similar levels of statistical significance. These results demonstrate that, even for flexible model specifications entailing high-dimensional covariates, the proposed method allows to produce

qualitatively similar results without extensive loss of significance, and we can thus confirm further robustness of the main empirical findings by Nunn and Wantchekon (2011).

7.2. Pre-Colonial Institutions and Regional Developments in Africa. Michalopoulos and Papaioannou (2013) analyze the effects of pre-colonial institutions on contemporary regional developments in Africa, controlling for various population, locational and geographic covariates including population density, distance to capital, distance to sea coast, distance to border, water area, land area, elevation, land suitable for agriculture, ecological suitability, petroleum, and diamond mine as well as country fixed effects in their baseline model. Estimates of these effects are obtained by running regressions with a sample that pools n_{ij} populated pixels $\ell \in [n_{ij}]$ in ethnic group i and country j across the cells $(i, j) \in [N] \times [M]$ of $N(= 93)$ ethnic groups and $M(= 48)$ countries. Standard errors are computed by the two-way cluster-robust method of Cameron, Gelbach and Miller (2011) for the ethnic group and district as two ways of clustering.

With our proposed method that is applicable to both high-dimensional models and multi-way clustering, they could consider even more flexible model specifications, for example, allowing for interactions of all combinations of geographical covariates and locational covariates. We present estimates with standard errors under such extended models with flexible specifications, demonstrate that qualitatively similar results continue to be obtained without substantial loss of statistical significance, and thus confirm further robustness of the main empirical findings by Michalopoulos and Papaioannou (2013).

Consider the model

$$Y_{ij,\ell} = D_{ij,\ell}\alpha + X'_{ij,\ell}\beta + R_{ij,\ell}^Y + \varepsilon_{ij,\ell}, \quad \mathbb{E}[\varepsilon_{ij,\ell} | D_{ij,\ell}, X_{ij,\ell}] = 0,$$

where $Y_{ij,\ell}$ denotes a regional development measured by night light intensity, $D_{ij,\ell}$ denotes an intensity measure of pre-colonial ethnic institutions, $X_{ij,\ell}$ contains population density, interactions of all combinations of locational controls (distance to capital, distance to sea coast, and distance to border), interactions of all combinations of geographical controls (water area, land area, elevation, land suitable for agriculture, ecological suitability, petroleum, and diamond mine), and country fixed effects, consisting of 82 or 83 dimensions of covariates in total. Note that the total number of regressors ($p + 1 = 83$ or 84) is much larger than the effective sample size ($\underline{C} = N \wedge M = 48$) of two-way clustering in this extended setting.

Table 3 summarize the estimates of the effects of pre-colonial institutions on regional development as measured by the “light density,” corresponding to parts of Table 3 of Michalopoulos and Papaioannou (2013). The last two columns in the table show the original estimates obtained under the prototypical model by Michalopoulos and Papaioannou (2013, Table 3) and corresponding lasso estimates obtained under more flexible model specification by our method. Across all the measures of pre-colonial institutions and all specifications, the original estimates and our lasso estimates are similar with similar levels of statistical significance. These results demonstrate that, even for flexible model specifications entailing high-dimensional covariates, the proposed method allows to produce qualitatively similar results without extensive loss of significance, and we can thus confirm further robustness of the main empirical findings by Michalopoulos and Papaioannou (2013).

8. CONCLUSION

In this paper, we investigate high-dimensional regression models when data is sampled under multi-way clustering. We establish the convergence rates for the lasso and post-lasso estimators under multi-way clustering. We then propose an inference method based on a post-double-selection procedure and show that it is asymptotically valid under multi-way clustering. Simulation studies demonstrate that the proposed procedure works well in finite sample under multi-way clustering in comparison with existing alternatives. We demonstrate that our method can enrich the flexibility of regression models and robustness of empirical results through a couple of empirical applications in growth and development economics.

Indeed, both multi-way clustering and high dimensionality are two important issues which concern applied research. The existing literature provide solutions to each of multi-way clustering and high-dimensionality separately. To our best knowledge, the literature does not seem to provide a solution to both of these issues simultaneously. In this paper, we filled this void in the literature.

APPENDIX A. MATHEMATICAL PROOFS

Throughout, the symbol $=$ stands for $\stackrel{a.s.}{=}$. We use the notations $Y = [Y_{11}, \dots, Y_{NM}]'$, $X = [X_{11}, \dots, X_{NM}]'$, $D = [D_{11}, \dots, D_{NM}]'$, $\mathcal{E} = [\varepsilon_{11}, \dots, \varepsilon_{NM}]'$, $V = [v_{11}, \dots, v_{NM}]'$, $R^Y = [R_{11}^Y, \dots, R_{NM}^Y]'$, $R^D = [R_{11}^D, \dots, R_{NM}^D]'$, $g = X\beta + R^Y$, and $m = X\gamma + R^D$. For any $A \subset [p]$, let $X_A = \{X_j : j \in A\}$, where X_j denotes the j -th the columns of X . Also define the projection operator by

$$\mathcal{P}_A = X_A(X_A'X_A)^-X_A',$$

and the orthogonal projection operator by $\mathcal{M}_A = I - \mathcal{P}_A$.

A.1. Proof of Lemma 1.

Proof. Our proof strategy closely follows that of Lemma D.2 in Davezies, D'Haultfoeuille and Guyonvarch (2018), except that we care about allowing the data generating process to vary with n to accommodate the increasing dimensionality p .

Under Assumption 1 (1) and (2), Lemma C.1 (a version of Aldous-Hoover representation) of Davezies, D'Haultfoeuille and Guyonvarch (2018) implies that, for each n , there exists a measurable function τ_n such that

$$\{Z_{ij}\}_{(i,j) \in \mathbb{N}^2} = \{\tau_n(U_{i0}, U_{0j}, U_{ij})\}_{(i,j) \in \mathbb{N}^2} \quad (\text{A.6})$$

holds, where $\{\{U_{i0}\}_{i \in \mathbb{N}}, \{U_{0j}\}_{j \in \mathbb{N}}, \{U_{ij}\}_{(i,j) \in \mathbb{N}^2}\}$ are i.i.d. uniform(0, 1) random variables.

The Hájek projection $H_n f$ of $\mathbb{G}_C f$ on the set \mathcal{G}_n is characterized by

$$\mathbb{E}\left[(\mathbb{G}_C f - H_n f) \cdot g(U_n)\right] = 0 \quad \text{for any } g(U_n) \in \mathcal{G}_n,$$

where $U_n = (U_{i0}, U_{0j})_{i \in [N], j \in [M]}$. Thus, for any U_c with $c = (c_1, c_2) \in \mathcal{I}_n = \{(i, 0), (0, j) : i \in [N], j \in [M]\}$, we have

$$\mathbb{E}[\mathbb{G}_C f | U_c] = \mathbb{E}[H_n f | U_c].$$

Because the range of H_n is a closed subspace, we have

$$H_n f = \sum_{i=1}^N \mathbb{E}[H_n f | U_{i0}] + \sum_{j=1}^M \mathbb{E}[H_n f | U_{0j}].$$

It follows from the above two equations that

$$H_n f = \sum_{i=1}^N \mathbb{E}[\mathbb{G}_C f | U_{i0}] + \sum_{j=1}^M \mathbb{E}[\mathbb{G}_C f | U_{0j}].$$

Now, fix $c \in \mathcal{I}_n$ and let $e(c) = (\mathbb{1}\{c_1 > 0\}, \mathbb{1}\{c_2 > 0\})$. By the independence of $\{\{U_{i0}\}_{i \in \mathbb{N}}, \{U_{0j}\}_{j \in \mathbb{N}}, \{U_{ij}\}_{(i,j) \in \mathbb{N}^2}\}$, Z_{ij} and U_c are independent whenever $c \neq (i, j) \odot e(c)$, where \odot denotes the Hadamard product. Thus, $\mathbb{E}[f(Z_{ij}) - \mathbb{E}f(Z_{11}) | U_c] = \mathbb{E}[f(Z_{ij}) - \mathbb{E}f(Z_{11})] = 0$. Therefore,

$$\begin{aligned} \mathbb{E}[\mathbb{G}_C f | U_c] &= \frac{\sqrt{C}}{NM} \sum_{i=1}^N \sum_{j=1}^M \mathbb{E}[f(Z_{ij}) - \mathbb{E}f(Z_{11}) | U_c] \\ &= \frac{\sqrt{C}}{NM} \sum_{i=1}^N \sum_{j=1}^M \mathbb{1}\{(i, j) \odot e(c) = c\} \mathbb{E}[f(Z_{ij}) - \mathbb{E}f(Z_{11}) | U_c]. \end{aligned}$$

The representation (A.6) implies, for all (i, j) such that $(i, j) \odot e(c) = c$,

$$\mathbb{E}[f(Z_{ij}) - \mathbb{E}f(Z_{11}) | U_c] = \mathbb{E}[f(Z_{c \vee 1}) - \mathbb{E}f(Z_{11}) | U_c],$$

i.e. the index outside the support of c can be changed to 1. Now, suppose $c_k = 0$, $k \in \{1, 2\}$.

The representation (A.6) again gives

$$\frac{\sqrt{C}}{NM} \sum_{i=1}^N \sum_{j=1}^M \mathbb{1}\{(i, j) \odot e(c) = c\} \mathbb{E}[f(Z_{c \vee 1}) - \mathbb{E}f(Z_{11}) | U_c] = \frac{\sqrt{C} C_k}{NM} \mathbb{E}[f(Z_{c \vee 1}) - \mathbb{E}f(Z_{11}) | U_c],$$

where $C_1 = N$, $C_2 = M$. Therefore,

$$\begin{aligned} H_n f &= \sum_{c \in \mathcal{I}_n} \frac{\sqrt{C} C_k}{NM} \mathbb{E}[f(Z_{c \vee 1}) - \mathbb{E}f(Z_{11}) | U_c] \\ &= \sum_{i=1}^N \frac{\sqrt{C}}{N} \mathbb{E}[f(Z_{i1}) - \mathbb{E}f(Z_{11}) | U_{i0}] + \sum_{j=1}^M \frac{\sqrt{C}}{M} \mathbb{E}[f(Z_{1j}) - \mathbb{E}f(Z_{11}) | U_{0j}], \end{aligned}$$

and each term in the two summands is independent from the others. This establishes the first claim of the lemma.

We next show that the variance of $H_n f$ can be calculated as

$$\begin{aligned} V(H_n f) &= \mu_N V(\mathbb{E}[f(Z_{11}) | U_{10}]) + \mu_M V(\mathbb{E}[f(Z_{11}) | U_{01}]) \\ &= \mu_N \text{Cov}(f(Z_{11}), f(Z_{12})) + \mu_M \text{Cov}(f(Z_{11}), f(Z_{21})). \end{aligned}$$

To see this, note that

$$\begin{aligned}
V(\mathbb{E}[f(Z_{11})|U_{10}]) &= \text{Cov}(\mathbb{E}[f(Z_{11})|U_{10}], \mathbb{E}[f(Z_{12})|U_{10}]) \\
&= \text{Cov}(f(Z_{11}), f(Z_{12})) - \mathbb{E}[\text{Cov}(f(Z_{11}), f(Z_{12})|U_{10})] \\
&= \text{Cov}(f(Z_{11}), f(Z_{12})) - \mathbb{E}[\text{Cov}\{f(\tau_n(U_{10}, U_{01}, U_{11})), f(\tau_n(U_{10}, U_{02}, U_{12}))|U_{10}\}] \\
&= \text{Cov}(f(Z_{11}), f(Z_{12})) - 0,
\end{aligned}$$

where the first equality follows from the representation (A.6), the second from the law of total covariance, the third from the representation (A.6), and the last from the fact that $\{\{U_{i0}\}_{i \in \mathbb{N}}, \{U_{0j}\}_{j \in \mathbb{N}}, \{U_{ij}\}_{(i,j) \in \mathbb{N}^2}\}$ are independent. Analogous lines of calculations yield $V(\mathbb{E}[f(X_{11})|U_{01}]) = \text{Cov}(f(X_{11}), f(X_{21}))$. Also, a direct calculation using Assumption 1 (1) and (2) shows

$$\begin{aligned}
V(\mathbb{G}_C f) &= \frac{\underline{C}}{(NM)^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} \text{Cov}(f(Z_{ij}), f(Z_{ij'})) + \frac{\underline{C}}{(NM)^2} \sum_{1 \leq i, i' \leq N} \sum_{j=1}^M \text{Cov}(f(Z_{ij}), f(Z_{i'j})) \\
&\quad - \frac{\underline{C}}{(NM)^2} \sum_{i=1}^N \sum_{j=1}^M V(f(Z_{ij})) \\
&= \mu_N \text{Cov}(f(Z_{11}), f(Z_{12})) + \mu_M \text{Cov}(f(Z_{11}), f(Z_{21})) + O\left(\frac{1}{\underline{C}}\right),
\end{aligned}$$

since $\underline{C}/NM \leq 1/\underline{C}$ and $\mathbb{E}f^2$ is bounded over n . This establishes the second claim of the lemma. \blacksquare

A.2. Proof of Theorem 1.

Proof. We will focusing on the result for $\tilde{\gamma}$ since the results for $\tilde{\eta}$ will follow analogously. The proof is divided into four steps. The conclusions from the first two steps give the rates for lasso. The third step provides bounds for the rates of the post-lasso in terms of convergence rates of lasso. The fourth step provides the ℓ_2 -norm rate.

Step 1. The oracle inequality follows directly from Lemma 6 of Belloni, Chen, Chernozhukov and Hansen (2012), which is applicable under Assumption 3 (1)–(2). This implies that we have the following bounds for lasso estimator:

$$\begin{aligned}\|\hat{\gamma} - \gamma\|_1 &\lesssim \frac{\sqrt{s}\lambda_2}{NM} + \frac{NM c_s^2}{\lambda_2} \quad \text{and} \\ \|X'_{ij}(\hat{\eta} - \eta)\|_n &\lesssim \frac{s\lambda_2}{NM} + c_s,\end{aligned}$$

conditionally on the regularized event $\lambda_2/NM \geq c\|(NM)^{-1} \sum_{i=1}^N \sum_{j=1}^M X_{ij}v_{ij}\|_\infty$ for some constant $c > 1$. Apply Assumption 3 (2) and use the choice of λ_2 to obtain the desired results in the first line.

Step 2. We now claim that, if we set $\lambda_2 = O\left(\sqrt{(NM)^2 \log a/\underline{C}}\right)$, then the regularized event

$$\max_{k \in [p]} \left| \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M X_{ij,k}v_{ij} - \mathbb{E}[X_{11,k}v_{11}] \right| \lesssim \frac{1}{c} \sqrt{\frac{\log a}{\underline{C}}} = \frac{\lambda_2}{NM}. \quad (\text{A.7})$$

realizes with probability at least $1 - C(\log \underline{C})^{-1}$.

First notice that the left-hand side can be bounded as

$$\begin{aligned}& \max_{k \in [p]} \left| \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M X_{ij,k}v_{ij} - \mathbb{E}[X_{11,k}v_{11}] \right| \\ & \leq \max_{k \in [p]} \left| \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M X_{ij,k}v_{ij} - \frac{M}{NM} \sum_{i=1}^N \mathbb{E}[X_{i1,k}v_{i1}|U_{i0}] - \frac{N}{NM} \sum_{j=1}^M \mathbb{E}[X_{1j,k}v_{1j}|U_{0j}] \right| \\ & \quad + \max_{k \in [p]} \left| \frac{M}{NM} \sum_{i=1}^N \mathbb{E}[X_{i1,k}v_{i1}|U_{i0}] - \mathbb{E}[X_{11,k}v_{11}] \right| + \max_{k \in [p]} \left| \frac{N}{NM} \sum_{j=1}^M \mathbb{E}[X_{1j,k}v_{1j}|U_{0j}] - \mathbb{E}[X_{11,k}v_{11}] \right| \\ & = (1) + (2) + (3)\end{aligned}$$

where $\mathbb{E}[X_{11,k}v_{11}] = 0$ is used. Part (1) is $O_P\left(\frac{1}{\underline{C}}\right)$ by Lemma 1 under Assumptions 1 and 2 (1). Using Lemma 2 (see Appendix B ahead), we can show (2) = $O_P\left(\sqrt{\frac{\log a}{N}}\right)$. To see this, note that Assumption 2 (1) implies that $\sigma^2 := \max_{k \in [p]} \frac{1}{N} \sum_{i=1}^N \mathbb{E}(\mathbb{E}[X_{i1}v_{i1}|U_{i0}])^2$ is

uniformly bounded, and Assumption 2 (1)–(2) suggests

$$\begin{aligned}
B^2 &= \mathbf{E}[\max_{i \in [N]} \max_{k \in [p]} (\mathbf{E}[X_{i1,k} v_{i1} | U_{i0}])^2] \\
&\leq \mathbf{E}[\max_{i \in [N]} \|X_{i1} v_{i1}\|_\infty^2] \\
&\leq \left(\mathbf{E}[\max_{i \in [N]} \|X_{i1} v_{i1}\|_\infty^q] \right)^{2/q} \\
&\leq N^{2/q} \left(\mathbf{E}[\|X_{11}\|_\infty^q |v_{11}|^q] \right)^{2/q} \\
&\leq N^{2/q} \left(\sqrt{\mathbf{E}[\|X_{11}\|_\infty^{2q}] \sqrt{\mathbf{E}[|v_{11}|^{2q}]} \right)^{2/q} = N^{2/q} B_n^2 O(1),
\end{aligned}$$

where the first inequality follows from the property of projection, the second inequality follows from Jensen's inequality, and the fourth inequality follows from Cauchy-Schwartz inequality. Thus, Lemma 2 (see Appendix B ahead) implies

$$|(2)| \lesssim_{\mathbf{P}} \sqrt{\frac{\log a}{N}} + \frac{B_n \log a}{N^{1-1/q}} \lesssim \sqrt{\frac{\log a}{N}}.$$

Finally, (3) = $O_{\mathbf{P}}(\sqrt{\frac{\log a}{M}})$ follows analogously. Therefore, (2) + (3) = $O_{\mathbf{P}}(\sqrt{\frac{\log a}{\underline{C}}})$.

Step 3. We now derive bounds for performance of post-lasso:

$$\|m_{ij} - X_{ij} \tilde{\gamma}\|_n \lesssim_{\mathbf{P}} \sqrt{\frac{s \log a}{\underline{C}}} + \frac{\|(I - \mathcal{P}_{\hat{I}_2})m\|}{\sqrt{NM}} \quad (\text{A.8})$$

where $m_{ij} = X'_{ij} \gamma + R_{ij}^D$. This part of proof closely follows the proof of Lemma 7 in Belloni, Chen, Chernozhukov and Hansen (2012) with some minor modifications. First note that $m - X_{\hat{I}_2} \tilde{\gamma} = (I - \mathcal{P}_{\hat{I}_2})m - \mathcal{P}_{\hat{I}_2} V$. This implies

$$\|m - X_{\hat{I}_2} \tilde{\gamma}\| \leq \|(I - \mathcal{P}_{\hat{I}_2})m\| + \|\mathcal{P}_{\hat{I}_2} V\|.$$

By the definition of $\mathcal{P}_{\hat{I}_2}$ and the operator norm,

$$\|\mathcal{P}_{\hat{I}_2} V\| \leq \|X_{\hat{I}_2} / \sqrt{NM} (X'_{\hat{I}_2} X_{\hat{I}_2} / NM)^{-1}\| \|X'_{\hat{I}_2} V / \sqrt{MM}\|$$

and

$$\|X_{\hat{I}_2} / \sqrt{NM} (X'_{\hat{I}_2} X_{\hat{I}_2} / NM)^{-1}\| \leq \sqrt{1 / \phi_{\min}(s + \tilde{m}_2)},$$

where $\tilde{m}_2 = |\hat{I}_2 \setminus T_2|$, $T_2 = \text{support}(\gamma)$. Thus under Assumption 4, we obtain

$$\|\mathcal{P}_{\hat{I}_2} V\| \leq \sqrt{1 / \phi_{\min}(s + \tilde{m}_2)} \|X'_{\hat{I}_2} V / \sqrt{MM}\| \leq \sqrt{\frac{s + \tilde{m}_2}{\phi_{\min}(s + \tilde{m}_2)}} \|XV / \sqrt{NM}\|_\infty \lesssim_{\mathbf{P}} \sqrt{\frac{s \log a}{\underline{C}}},$$

where the last inequality follows from equation (A.7) and Lemma 3 (see Appendix B ahead). This shows (A.8).

By (A.7), the same argument as that of Lemma 7 in Belloni, Chen, Chernozhukov and Hansen (2012) establishes

$$\frac{\|(I - \mathcal{P}_{\hat{I}_2})m\|}{\sqrt{NM}} \lesssim \frac{\sqrt{s}\lambda_2}{NM} + c_s.$$

Therefore, (A.8) can be rewritten as

$$\|m_{ij} - X_{ij}\tilde{\gamma}\|_n \lesssim_P \sqrt{\frac{s \log a}{\underline{C}}} + \frac{\sqrt{s}\lambda_2}{NM} + c_s.$$

Next, applying Lemma 3 (see Appendix B ahead), we have

$$\|\tilde{\gamma} - \gamma\|_1 \leq \sqrt{\|\tilde{\gamma} - \gamma\|_0} \|\tilde{\gamma} - \gamma\| \leq \sqrt{s + \tilde{m}_2} \frac{\|X_{ij}(\tilde{\gamma} - \gamma)\|_n}{\sqrt{\phi_{\min}(s + \tilde{m}_2)}}$$

under Assumption 3 (1) and 4.

Combining the above bounds, the choice of λ_2 , and Assumption 3 (2), we obtain

$$\begin{aligned} \|\tilde{\gamma} - \gamma\|_1 &\lesssim_P \sqrt{\frac{s^2 \log a}{\underline{C}}} \quad \text{and} \\ \|X_{ij}(\tilde{\gamma} - \gamma)\|_n &\lesssim_P \sqrt{\frac{s \log a}{\underline{C}}}. \end{aligned}$$

Step 4. The ℓ_2 -norm rates are implied by the prediction norm rates, Assumption 4, and Lemma 3 (see Appendix B ahead). ■

A.3. Proof of Theorem 2.

Proof. Our proof follows parallel steps to Steps 1-6 in the proof of Theorem 1 in Belloni, Chernozhukov and Hansen (2014). However, due to the two-way cluster sampling, most of the probabilistic bounds are established differently.

We use the short-hand notation

$$\tilde{b}_Z(A) := \operatorname{argmin}_{b \in \mathbb{R}^p: b_j = 0 \forall j \in A^c} \|Z - X'b\|^2$$

for any vector $Z \in \mathbb{R}^n$.

Step 1 Write $\tilde{\alpha} = [D' \mathcal{M}_{\hat{\tau}} D / NM]^{-1} D' \mathcal{M}_{\hat{\tau}} Y / NM$ and thus we can write

$$\sqrt{\underline{C}}(\tilde{\alpha} - \alpha) = [D' \mathcal{M}_{\hat{\tau}} D / NM]^{-1} \cdot \sqrt{\underline{C}} D' \mathcal{M}_{\hat{\tau}} (g + \mathcal{E}) / NM =: (II)^{-1} \cdot (I).$$

By Steps 2 and 3 to be presented below, we obtain

$$(II) = V'V / NM + o_P(1) \text{ and } (I) = \sqrt{\underline{C}} V' \mathcal{E} / NM + o_P(1).$$

Also note that $V'V / NM = E[v_{11}^2] + o_P(1)$ by Lemma 1 and Assumption 1–2, which can be shown following the same arguments as those in Step 3 of the proof for Theorem 3. Under Assumption 2 (1), $E[v_{11}^2]$ is bounded and bounded away from zero uniformly in n . Therefore $(II)^{-1} = E[v_{11}^2]^{-1} + o_P(1)$.

Under Assumption 2 (3), σ^2 is bounded and bounded away from zero. Setting $W_{ij} := \sigma^{-1} v_{ij} \varepsilon_{ij}$ and $Z_{ij} \xrightarrow{f} W_{ij}$, we have $E f(Z_{11}) = 0$ and

$$\mathbb{G}_C f = \frac{\sqrt{\underline{C}}}{NM} \sum_{i=1}^N \sum_{j=1}^M W_{ij} = \sigma^{-1} \sqrt{\underline{C}} (\tilde{\alpha} - \alpha) + o_P(1).$$

Apply Lemma 1 under Assumption 1 and 2 (1) to obtain the Hájek projection

$$H_n f = \sum_{i=1}^N \frac{\sqrt{\underline{C}}}{N} E[f(Z_{i1}) | U_{i0}] + \sum_{j=1}^M \frac{\sqrt{\underline{C}}}{M} E[f(Z_{1j}) | U_{0j}]$$

of $\mathbb{G}_C f$, where terms in each summand are independent and two summands are independent of each other. We now check Lyapunov's conditions. First, note that Assumption 2 (1) guarantees that the third moments of both summands are bounded uniformly in n . Second, the second part of Lemma 1 implies that

$$\begin{aligned} \lim_{n \rightarrow \infty} V(H_n f) &= \bar{\mu}_N V(E[f(Z_{11}) | U_{10}]) + \bar{\mu}_M V(E[f(Z_{11}) | U_{01}]) \\ &= \bar{\mu}_N E[f(Z_{11}) f(Z_{12})] + \bar{\mu}_M E[f(Z_{11}) f(Z_{21})] = \Gamma \in (c, \infty) \end{aligned}$$

a.s. for $c > 0$, where the last inequalities follow from Assumption 2 (3). Therefore, we apply Lyapunov's CLT to obtain

$$H_n f \rightsquigarrow N\left(0, \bar{\mu}_N V(E[f(Z_{11}) | U_{10}]) + \bar{\mu}_M V(E[f(Z_{11}) | U_{01}])\right).$$

The first equality in the variance equation of Lemma 1 yields

$$V(\mathbb{G}_C f) = \bar{\mu}_N E[f(Z_{11}) f(Z_{12})] + \bar{\mu}_M E[f(Z_{11}) f(Z_{21})] + o_P(1),$$

where the right-hand side is asymptotically positive and bounded away from zero. Therefore,

$$\sqrt{\underline{C}}(\tilde{\alpha} - \alpha) = \mathbb{G}_C f \rightsquigarrow N(0, \sigma^2).$$

Step 2 Use $D = m + V$ to decompose

$$\begin{aligned} (I) &= \sqrt{\underline{C}}V'\mathcal{E}/NM + \sqrt{\underline{C}}m'\mathcal{M}_{\hat{\Gamma}}g/NM + \sqrt{\underline{C}}m'\mathcal{M}_{\hat{\Gamma}}\mathcal{E}/NM + \sqrt{\underline{C}}V'\mathcal{M}_{\hat{\Gamma}}g/NM - \sqrt{\underline{C}}V'\mathcal{P}_{\hat{\Gamma}}\mathcal{E}/NM \\ &= \sqrt{\underline{C}}V'\mathcal{E}/NM + (1a) + (1b) + (1c) - (1d). \end{aligned}$$

By Steps 5 and 6 to be presented below, we have

$$|(1a)| \lesssim \sqrt{\underline{C}}\|\mathcal{M}_{\hat{\Gamma}}m/\sqrt{NM}\| \cdot \|\mathcal{M}_{\hat{\Gamma}}g/\sqrt{NM}\| \lesssim_{\mathbb{P}} \sqrt{\frac{s^2(\log a)^2}{\underline{C}}}.$$

Using the decompositions $m = X\gamma + R^D$, $m'\mathcal{P}_{\hat{\Gamma}} = \tilde{b}'_m(\hat{I})X'$ and $m'\mathcal{M}_{\hat{\Gamma}}\mathcal{E} = (R^D)'\mathcal{E} - (\tilde{b}_m(\hat{I}) - \gamma)'X'\mathcal{E}$, one has

$$|(1b)| \leq \sqrt{\underline{C}}|R^{D'}\mathcal{E}/NM| + \sqrt{\underline{C}}|(\tilde{b}_m(\hat{I}) - \gamma)X'\mathcal{E}/NM| \lesssim_{\mathbb{P}} \sqrt{\frac{s^2(\log a)^2}{\underline{C}}},$$

because, under Assumptions 2 (1) and 3 (2),

$$\sqrt{\underline{C}}|R^{D'}\mathcal{E}/NM| \leq \sqrt{\underline{C}}\sqrt{R^{D'}R^D/NM} \cdot O_{\mathbb{P}}\left(\sqrt{\frac{1}{NM}\mathbb{E}\|\mathcal{E}\|^2}\right) \lesssim_{\mathbb{P}} \sqrt{\frac{s}{\underline{C}}}$$

and

$$\begin{aligned} \sqrt{\underline{C}}|(\tilde{b}_m(\hat{I}) - \gamma)X'\mathcal{E}/NM| &\leq \sqrt{\underline{C}}\|\tilde{b}_m(\hat{I}) - \gamma\|_1\|X'\mathcal{E}/NM\|_{\infty} \\ &\lesssim_{\mathbb{P}} \sqrt{\underline{C}}\sqrt{\frac{s^2 \log a}{\underline{C}}} \cdot \sqrt{\frac{\log a}{\underline{C}}} = \sqrt{\frac{s^2(\log a)^2}{\underline{C}}}, \end{aligned}$$

where $\|\tilde{b}_m(\hat{I}) - \gamma\|_1 \lesssim_{\mathbb{P}} \sqrt{\frac{s^2 \log a}{\underline{C}}}$ follows from Step 5 and $\|X'\mathcal{E}/NM\|_{\infty} \lesssim_{\mathbb{P}} \sqrt{\frac{\log a}{\underline{C}}}$ follows from Step 4. Third, using the same argument as above for (1b) following Steps 4 and 6 and $g = X\beta + R^Y$, we have

$$|(1c)| \leq \sqrt{\underline{C}}|R^{Y'}V/\sqrt{NM}| + \sqrt{\underline{C}}|(\tilde{b}_g(\hat{I}) - \beta)'X'V/\sqrt{NM}| \lesssim_{\mathbb{P}} \sqrt{\frac{s}{\underline{C}}} + \sqrt{\frac{s^2(\log a)^2}{\underline{C}}}.$$

Finally,

$$|(1d)| \leq \sqrt{\underline{C}}|\tilde{b}_V(\hat{I})'X'\mathcal{E}/NM| \leq \sqrt{\underline{C}}\|\tilde{b}_V(\hat{I})\|_1\|X'\mathcal{E}/NM\|_{\infty} \lesssim_{\mathbb{P}} \sqrt{\frac{s^2(\log a)^2}{\underline{C}}}$$

following equation (A.7) in the proof of Theorem 1, and

$$\begin{aligned} \|\tilde{b}_V(\hat{I})\|_1 &\leq \sqrt{\hat{s}}\|\tilde{b}_V(\hat{I})\| \leq \sqrt{\hat{s}}\|(X'_{\hat{I}}X_{\hat{I}}/NM)^{-1}X'_{\hat{I}}V/NM\| \\ &\lesssim_{\mathbb{P}} \frac{\sqrt{\hat{s}}}{\phi_{\min}(C\hat{s})}\sqrt{\hat{s}}\|X'_{\hat{I}}V/NM\|_{\infty} \lesssim_{\mathbb{P}} \sqrt{\frac{s^2 \log a}{\underline{C}}} \end{aligned}$$

following Step 4, Lemma 3, and Assumption 4.

Step 3 We can write

$$\begin{aligned} (II) &= (m+V)' \mathcal{M}_{\hat{I}}(m+V)/NM \\ &= V'V/NM + m' \mathcal{M}_{\hat{I}}m/NM + 2m' \mathcal{M}_{\hat{I}}V/NM - V' \mathcal{P}_{\hat{I}}V/NM \\ &= V'V/NM + (2a) + (2b) - (2c). \end{aligned}$$

We have $|(2a)| \lesssim_{\mathbb{P}} \sqrt{\frac{s^2(\log a)^2}{\underline{C}}}$ by Step 5, $|(2b)| \lesssim_{\mathbb{P}} \sqrt{\frac{s^2(\log a)^2}{\underline{C}}}$ by a similar argument to bounding $|(1b)|$, and $|(2c)| \lesssim_{\mathbb{P}} \sqrt{\frac{s^2(\log a)^2}{\underline{C}}}$ by a similar argument to bounding $|(1d)|$.

Step 4 In this step, we show that the following regularized events hold with probability $1 - o(1)$:

$$(a) \sqrt{\underline{C}}\|X'\mathcal{E}/NM\|_{\infty} \lesssim_{\mathbb{P}} \sqrt{\log a} \text{ and } (b) \sqrt{\underline{C}}\|X'V/NM\|_{\infty} \lesssim_{\mathbb{P}} \sqrt{\log a}.$$

This claim follows from similar lines of argument to those showing equation (A.7) in the proof of Theorem 1 under Assumptions 1 and 2 (1)–(2).

Step 5 In this step, we show

$$(a) \|\mathcal{M}_{\hat{I}}m/\sqrt{NM}\| \lesssim_{\mathbb{P}} \sqrt{\frac{s \log a}{\underline{C}}} \text{ and } (b) \|\tilde{b}_m(\hat{I}) - \gamma\| \lesssim_{\mathbb{P}} \sqrt{\frac{s \log a}{\underline{C}}}.$$

First, by applying Theorem 1 under Assumptions 1, 2 (1)–(2), 3(1)–(2), and 4, and by following the same argument as the one in Step 5 of Belloni, Chernozhukov and Hansen (2014), we have

$$\|\mathcal{M}_{\hat{I}}m/\sqrt{NM}\| \leq \|\mathcal{M}_{\hat{I}_2}m/\sqrt{NM}\| \leq \|(X\tilde{b}_D(\hat{I}_2) - m)/\sqrt{NM}\| \lesssim_{\mathbb{P}} \sqrt{\frac{s \log a}{\underline{C}}},$$

where the first inequality follows from $\hat{I}_2 \subset \hat{I}$ and the second follows from the fact that $\tilde{b}_m(\hat{I}_2)$ minimizes $\|m - X'_{\hat{I}_2}b\|$. Second, the reverse triangle inequality yields

$$\left| \|X(\tilde{b}_m(\hat{I}) - \gamma)/\sqrt{NM}\| - \|R_m/\sqrt{NM}\| \right| \lesssim_{\mathbb{P}} \|\mathcal{M}_{\hat{I}}m/\sqrt{NM}\|,$$

and, by Assumption 3 (2), $\|R_m/\sqrt{NM}\| \lesssim_P \sqrt{s/\underline{C}}$. Thus, by using Lemma 3 with Assumptions 3 (1) and 4, we obtain

$$\begin{aligned} \|\tilde{b}_m(\hat{I}) - \gamma\| &\lesssim_P \sqrt{\phi_{\min}(\hat{s} + s)} \|\tilde{b}_m(\hat{I}) - \gamma\| \\ &\leq \|X(\tilde{b}_m(\hat{I}) - \gamma)/\sqrt{NM}\| \lesssim_P \sqrt{\frac{s \log a}{\underline{C}}}. \end{aligned}$$

Step 6 Finally, we can show

$$(a) \sqrt{\underline{C}} \|\mathcal{M}_{\hat{T}} g / NM\| \lesssim_P \sqrt{\frac{s \log a}{\underline{C}}} \text{ and } (b) \|\tilde{b}_g(\hat{I}) - \beta\| \lesssim_P \sqrt{\frac{s \log a}{\underline{C}}}.$$

following similar lines of argument to those of Step 5 under Assumptions 1, 2, 3, and 4. ■

A.4. Proof of Theorem 3.

Proof. First, note that we have the following decomposition

$$\begin{aligned} &|\hat{Q}^{-1} \hat{\Gamma} \hat{Q}^{-1} - Q^{-1} \Gamma Q^{-1}| \\ &\lesssim |\hat{Q}^{-1} - Q^{-1}| |\hat{Q}^{-1} + Q^{-1}| |\hat{\Gamma}| + |\hat{\Gamma} - \Gamma| |Q^{-1}|^2, \end{aligned}$$

where $|Q^{-1}|$ is bounded away from zero uniformly by Assumption 2(1). The rest of this proof is divide into 5 steps. In Steps 1 and 2, we obtain a bound for $|\hat{\Gamma} - \Gamma|$. In Steps 3 and 4, we obtain a bound for $|\hat{Q}^{-1} - Q^{-1}|$. Finally, Step 5 shows a bound for $|\hat{\Gamma}|$ and $|\hat{Q}^{-1} + Q^{-1}|$.

Step 1. We derive a bound for $|\tilde{\Gamma} - \Gamma|$, where

$$\tilde{\Gamma} = \frac{\mu_N}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} v_{ij} \varepsilon_{ij} v_{ij'} \varepsilon_{ij'} + \frac{\mu_M}{N^2 M} \sum_{1 \leq i, i' \leq N} \sum_{j=1}^M v_{ij} \varepsilon_{ij} v_{i'j} \varepsilon_{i'j}.$$

We first claim that

$$\frac{\mu_N}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} v_{ij} \varepsilon_{ij} v_{ij'} \varepsilon_{ij'} = \bar{\mu}_N \mathbf{E}[v_{11} \varepsilon_{11} v_{12} \varepsilon_{12}] + o_P(1). \quad (\text{A.9})$$

Note that, for each n , for any $i, \iota \in [N]$ and $j, k, l, m \in [M]$, we have

$$\begin{aligned} \text{Cov}\left(v_{ij} \varepsilon_{ij} v_{ik} \varepsilon_{ik}, v_{\iota l} \varepsilon_{\iota l} v_{\iota m} \varepsilon_{\iota m}\right) &\leq \max_{i \in [N], j, j' \in [M]} V(v_{ij} \varepsilon_{ij} v_{ij'} \varepsilon_{ij'}) \\ &= \max_{i \in [N], j, j' \in [M]} \left\{ \mathbf{E}[(v_{ij} \varepsilon_{ij} v_{ij'} \varepsilon_{ij'})^2] - (\mathbf{E}[v_{ij} \varepsilon_{ij} v_{ij'} \varepsilon_{ij'}])^2 \right\}. \end{aligned}$$

Using Cauchy-Schwartz's inequality with Assumptions 1 (1) and 2 (1), the first term in the variance can be bounded as

$$\begin{aligned} \mathbb{E}[(v_{ij}\varepsilon_{ij}v_{ij'}\varepsilon_{ij'})^2] &\leq \sqrt{\mathbb{E}[v_{ij}^4\varepsilon_{ij}^4]\mathbb{E}[v_{ij'}^4\varepsilon_{ij'}^4]} \\ &\leq \sqrt{\mathbb{E}v_{11}^8\mathbb{E}\varepsilon_{11}^8} = O(1) \end{aligned}$$

uniformly over n . Under Assumptions 1 (1) and 2 (1), similar calculations can be carried out to the square-root of the second term to obtain

$$\mathbb{E}[v_{ij}\varepsilon_{ij}v_{ij'}\varepsilon_{ij'}] \leq \sqrt{\mathbb{E}[v_{ij}^2\varepsilon_{ij}^2]\mathbb{E}[v_{ij'}^2\varepsilon_{ij'}^2]} \leq \sqrt{\mathbb{E}[v_{11}^4]\mathbb{E}[\varepsilon_{11}^4]} = O(1)$$

uniformly over n . This shows that, for any n , for any $i, \iota \in [N]$ and $j, k, l, m \in [M]$, it holds that, for a $K > 0$ independent of n ,

$$\left| \text{Cov}\left(v_{ij}\varepsilon_{ij}v_{ik}\varepsilon_{ik}, v_{il}\varepsilon_{il}v_{im}\varepsilon_{im}\right) \right| \leq K. \quad (\text{A.10})$$

With this bound of the covariance, we can bound the variance as

$$\begin{aligned} &V\left(\frac{1}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} v_{ij}\varepsilon_{ij}v_{ij'}\varepsilon_{ij'}\right) \\ &= \text{Cov}\left(\frac{1}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} v_{ij}\varepsilon_{ij}v_{ij'}\varepsilon_{ij'}, \frac{1}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} v_{ij}\varepsilon_{ij}v_{ij'}\varepsilon_{ij'}\right) \\ &= \frac{1}{N^2M^4} \sum_{i=1}^N \sum_{1 \leq j, k, l, m \leq M} \text{Cov}\left(v_{ij}\varepsilon_{ij}v_{ik}\varepsilon_{ik}, v_{il}\varepsilon_{il}v_{im}\varepsilon_{im}\right) \\ &\quad + \frac{2}{N^2M^4} \sum_{j=1}^M \sum_{1 \leq i, i' \leq N} \sum_{1 \leq k, l \leq M} \text{Cov}\left(v_{ij}\varepsilon_{ij}v_{ik}\varepsilon_{ik}, v_{i'j}\varepsilon_{i'j}v_{i'l}\varepsilon_{i'l}\right) + o\left(\frac{1}{C}\right) \\ &= O\left(\frac{1}{C}\right) = o(1) \end{aligned}$$

uniformly over n , where the second equality follows from Assumption 1 (2) and counting the number of terms in each summand, and the third equality is due to (A.10). Applying Chebyshev's inequality, it follows that

$$\begin{aligned} &\mathbb{P}\left(\left|\frac{1}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} (v_{ij}\varepsilon_{ij}v_{ij'}\varepsilon_{ij'} - \mathbb{E}[v_{ij}\varepsilon_{ij}v_{ij'}\varepsilon_{ij'}])\right| > \epsilon\right) \\ &\leq \frac{\sup_n V\left(\frac{1}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} v_{ij}\varepsilon_{ij}v_{ij'}\varepsilon_{ij'}\right)}{\epsilon^2} = \frac{1}{\epsilon^2} \cdot o(1). \quad (\text{A.11}) \end{aligned}$$

for any $\epsilon > 0$.

Also, under Assumption 2 (1), the first result in Lemma D.10 of Davezies, D'Haultfoeuille and Guyonvarch (2018) ensures

$$\mathbb{E} \left[\left| \frac{\mu_N}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} v_{ij} \varepsilon_{ij} v_{ij'} \varepsilon_{ij'} - \frac{\mu_N}{NM(M-1)} \sum_{i=1}^N \sum_{j=1}^M \sum_{j' \neq j} v_{ij} \varepsilon_{ij} v_{ij'} \varepsilon_{ij'} \right| \right] = o(1)$$

uniformly over n . Furthermore, under Assumptions 1 (1) and 2 (1), we have

$$\mathbb{E} \left[\frac{1}{NM(M-1)} \sum_{i=1}^N \sum_{j=1}^M \sum_{j' \neq j} v_{ij} \varepsilon_{ij} v_{ij'} \varepsilon_{ij'} \right] = \mathbb{E}[v_{11} \varepsilon_{11} v_{12} \varepsilon_{12}].$$

Combining these with (A.11), we obtain (A.9). A symmetric argument also shows

$$\frac{\mu_M}{N^2 M} \sum_{1 \leq i, i' \leq N} \sum_{j=1}^M v_{ij} \varepsilon_{ij} v_{i'j} \varepsilon_{i'j} = \bar{\mu}_M \mathbb{E}[v_{11} \varepsilon_{11} v_{21} \varepsilon_{21}] + o_P(1).$$

Therefore, we obtain $|\tilde{\Gamma} - \Gamma| = o_P(1)$.

Step 2. In this step we bound $|\hat{\Gamma} - \tilde{\Gamma}|$, where $\tilde{\Gamma}$ is defined in Step 1. Consider the decomposition

$$\underbrace{\frac{\mu_N}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} \left(\hat{v}_{ij} \hat{\varepsilon}_{ij} \hat{v}_{ij'} \hat{\varepsilon}_{ij'} - v_{ij} \varepsilon_{ij} v_{ij'} \varepsilon_{ij'} \right)}_{(1)} + \underbrace{\frac{\mu_N}{N^2 M} \sum_{1 \leq i, i' \leq N} \sum_{j=1}^M \left(\hat{v}_{ij} \hat{\varepsilon}_{ij} \hat{v}_{i'j} \hat{\varepsilon}_{i'j} - v_{ij} \varepsilon_{ij} v_{i'j} \varepsilon_{i'j} \right)}_{(2)}.$$

Recall $\varepsilon_{ij} = Y_{ij} - Z'_{ij} \eta - R_{ij}^Y$, $\hat{\varepsilon}_{ij} = Y_{ij} - Z'_{ij} \tilde{\eta} - R_{ij}^Y$, $v_{ij} = D_{ij} - X'_{ij} \gamma - R_{ij}^D$, $\hat{v}_{ij} = D_{ij} - X'_{ij} \tilde{\gamma} - R_{ij}^D$, and thus, $\hat{\varepsilon}_{ij} - \varepsilon_{ij} = Z'_{ij} (\tilde{\eta} - \eta) - R_{ij}^Y$ and $\hat{v}_{ij} - v_{ij} = X'_{ij} (\tilde{\gamma} - \gamma) - R_{ij}^D$. We can further decompose (1) as

$$\begin{aligned} (1) &= \frac{\mu_N}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} \left((\hat{v}_{ij} - v_{ij}) \hat{\varepsilon}_{ij} \hat{v}_{ij'} \hat{\varepsilon}_{ij'} \right) + \frac{\mu_N}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} \left(v_{ij} (\hat{\varepsilon}_{ij} - \varepsilon_{ij}) \hat{v}_{ij'} \hat{\varepsilon}_{ij'} \right) \\ &\quad + \frac{\mu_N}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} \left(v_{ij} \varepsilon_{ij} (\hat{v}_{ij'} - v_{ij'}) \hat{\varepsilon}_{ij'} \right) + \frac{\mu_N}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} \left(v_{ij} \varepsilon_{ij} v_{ij'} (\hat{\varepsilon}_{ij'} - \varepsilon_{ij'}) \right) \\ &= (1a) + (1b) + (1c) + (1d). \end{aligned}$$

Under Assumption 2 (1), we first bound

$$\begin{aligned}
(1a) &\lesssim_{\mathbb{P}} \left| \frac{\mu_N}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} \left((\widehat{v}_{ij} - v_{ij}) \widehat{\varepsilon}_{ij} \widehat{v}_{ij'} \widehat{\varepsilon}_{ij'} \right) \right| \\
&\leq \left| \frac{\mu_N}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} X'_{ij} (\widehat{\gamma} - \gamma) \widehat{\varepsilon}_{ij} \widehat{v}_{ij'} \widehat{\varepsilon}_{ij'} \right| + \left| \frac{\mu_N}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} R_{ij}^D \widehat{\varepsilon}_{ij} \widehat{v}_{ij'} \widehat{\varepsilon}_{ij'} \right| \\
&=(1aa) + (1ab).
\end{aligned}$$

We obtain

$$\begin{aligned}
(1aa) &= \left| \frac{\mu_N}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} X'_{ij} (\widehat{\gamma} - \gamma) \widehat{\varepsilon}_{ij} \widehat{v}_{ij'} \widehat{\varepsilon}_{ij'} \right| \\
&\leq \left| \frac{\mu_N}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} X'_{ij} (\widehat{\gamma} - \gamma) \cdot \left(Z'_{ij} (\widehat{\eta} - \eta) + \varepsilon_{ij} - R_{ij}^Y \right) \right. \\
&\quad \left. \cdot \left(X'_{ij'} (\widehat{\gamma} - \gamma) + v_{ij'} - R_{ij'}^D \right) \cdot \left(Z'_{ij'} (\widehat{\eta} - \eta) + \varepsilon_{ij'} - R_{ij'}^Y \right) \right| = o_{\mathbb{P}}(1),
\end{aligned}$$

where the last equality follows from triangle inequality, Cauchy-Schwartz's inequality, Theorem 1, Assumptions 2 (1)–(2) and 3 (2), and the rate conditions in the statement of the theorem. To see this, note that, under Assumption 2 (2), Theorem 1, and the rate condition in the theorem, we have

$$\begin{aligned}
&\sqrt{\frac{\underline{C}}{N^2 M^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} \left(X'_{ij} (\widehat{\gamma} - \gamma) Z'_{ij} (\widehat{\eta} - \eta) \right)^2} \\
&\leq \sqrt{\frac{\underline{C}}{NM} \max_{i \in [N], j \in [M]} \|Z_{ij}\|_{\infty}^2 \|\widehat{\eta} - \eta\|_1^2 M \|X'_{ij} (\widehat{\gamma} - \gamma)\|_n^2} \\
&\lesssim_{\mathbb{P}} \sqrt{\frac{(NM)^{1/q} B_n^2 s^3 (\log a)^2}{\underline{C} N}} = o(1).
\end{aligned}$$

Furthermore, by Theorem 1 and Assumptions 2 (1) and 3 (2), we have

$$\begin{aligned} & \sqrt{\frac{\underline{C}}{N^2 M^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} \left(X'_{ij} (\hat{\gamma} - \gamma) \varepsilon_{ij} \right)^2} \\ & \leq \sqrt{\frac{\underline{C}}{NM} \max_{i \in [N], j \in [M]} |\varepsilon_{ij}|^2 M \|X'_{ij} (\hat{\gamma} - \gamma)\|_n^2} \\ & \lesssim_P \sqrt{\frac{(NM)^{1/q_s} \log a}{N}} = o(1). \end{aligned}$$

The rest of the terms can be shown to be of smaller orders using similar arguments. Finally, the rate condition from the statement of the theorem gives

$$\sqrt{\frac{\underline{C}}{N^2 M^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} \left(R_{ij}^D R_{ij}^Y \right)^2} \leq \sqrt{\|R_{ij}^D R_{ij}^Y\|_n^2} = O(1).$$

Similarly, using Cauchy-Schwartz's inequality, Theorem 1, Assumptions 2 (1)–(2) and 3 (2), and the addition rate conditions in the statement of the theorem, we obtain

$$(1ab) = \left| \frac{\mu N}{NM^2} \sum_{i=1}^N \sum_{1 \leq j, j' \leq M} R_{ij}^D \hat{\varepsilon}_{ij} \hat{v}_{ij'} \hat{\varepsilon}_{ij'} \right| = o_P(1).$$

These results yield (1a) = $o_P(1)$. Following analogous but simpler arguments, we can show that (1b), (1c) and (1d) are $o_P(1)$. This shows (1) = $o_P(1)$. Similar lines of argument under the same set of assumptions show (2) = $o_P(1)$.

Step 3. In this and the next steps, we bound $|\hat{Q}^{-1} - Q^{-1}|$. Note that

$$|\hat{Q} - Q| = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M \hat{v}_{ij}^2 - \mathbb{E}[v_{11}^2] = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M (v_{ij}^2 - \mathbb{E}[v_{11}^2]) + \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M (\hat{v}_{ij}^2 - v_{ij}^2). \quad (\text{A.12})$$

The current step bounds the first term on the right-hand side, and Step 4 below bounds the second term on the right-hand side. With the notation $f(Z_{ij}) = v_{ij}^2$, the first term on right-hand side becomes

$$\frac{1}{\sqrt{\underline{C}}} \mathbb{G}_C f = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M (v_{ij}^2 - \mathbb{E}[v_{11}^2]).$$

Applying Lemma 1 under Assumptions 1 and 2 (1) suggests that its Hájek projection equals

$$\begin{aligned} \frac{1}{\sqrt{\underline{C}}} H_n f &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}[v_{i1}^2 - \mathbb{E}v_{i1}^2 | U_{i0}] + \frac{1}{M} \sum_{i=1}^N \mathbb{E}[v_{1j}^2 - \mathbb{E}v_{1j}^2 | U_{0j}] \\ &= O_{\mathbb{P}}\left(\frac{1}{\sqrt{N}} + \frac{1}{\sqrt{M}}\right) = O_{\mathbb{P}}\left(\frac{1}{\sqrt{\underline{C}}}\right), \end{aligned}$$

where the second equality follows from Lyapunov's CLT applied under Assumption 2 (1) – note that Assumption 2 implies the third moments for both terms of the right-hand side to be bounded. Assumption 2 and the second claim in Lemma 1 imply that $V(H_n f) = V(\mathbb{G}_C f) + O(\underline{C}^{-1})$. Since $H_n f$ is a projection of $\mathbb{G}_C f$, we obtain $\frac{1}{\sqrt{\underline{C}}}\mathbb{G}_C f = O_{\mathbb{P}}(\underline{C}^{-1/2}) = o_{\mathbb{P}}(1)$.

Step 4. To bound the second term on the RHS of equation (A.12), note that

$$\begin{aligned} \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M (\hat{v}_{ij}^2 - v_{ij}^2) &= \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M [(X_{ij}\gamma)^2 - (X_{ij}\hat{\gamma})^2] + \frac{2}{NM} \sum_{i=1}^N \sum_{j=1}^M D_{ij} X'_{ij} (\hat{\gamma} - \gamma) \\ &\quad + \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M (R_{ij}^D)^2 - \frac{2}{NM} \sum_{i=1}^N \sum_{j=1}^M D_{ij} R_{ij}^D + \frac{2}{NM} \sum_{i=1}^N \sum_{j=1}^M R_{ij}^D X'_{ij} \gamma \\ &= (3a) + (3b) + (3c) + (3d) + (3e). \end{aligned}$$

The first term can be bounded by

$$\begin{aligned} |(3a)| &= \left| \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M [(X_{ij}\gamma)^2 - (X_{ij}\hat{\gamma})^2] \right| \\ &\lesssim_{\mathbb{P}} \sup_{\substack{\|\delta\|=1 \\ \|\delta\|_0 \leq Cs}} \delta' \left(\frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M X_{ij} X'_{ij} \right) \delta \cdot \|\hat{\gamma} - \gamma\|^2 + 2 \sqrt{\frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M (X'_{ij}\gamma)^2} \cdot \|\hat{\gamma} - \gamma\| \\ &\leq \sqrt{\phi_{\max}(Cs)} \|\hat{\gamma} - \gamma\|^2 + 2 O_{\mathbb{P}} \left(\sqrt{\frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M \mathbb{E}(X'_{ij}\gamma)^2} \right) \|\hat{\gamma} - \gamma\| \lesssim_{\mathbb{P}} \sqrt{\frac{s \log a}{\underline{C}}}, \end{aligned}$$

where the first inequality follows from Assumption 3 (1), Lemma 3, and Lemma 3.1 in the supplementary appendix of van de Geer, Bühlmann, Ritov and Dezeure (2014), and the third follows from Theorem 1 and Assumption 4. An application of Cauchy-Schwartz's inequality and Theorem 1 gives $(3b) = O_{\mathbb{P}}(\sqrt{s \log a / \underline{C}})$. $(3c) = s / \underline{C}$ follows from Assumption 3. Cauchy-Schwartz's inequality and Assumptions 2(1), 3(2) lead to $(3d) \leq \|(NM)^{-1/2} R^D\| \sqrt{\mathbb{E}[D_{11}^2]} = O_{\mathbb{P}}(\sqrt{s / \underline{C}}) O(1)$. Using the property of the projection and a

similar argument to that of (3d), we conclude that $(3e) \leq \|(NM)^{-1/2}R^D\|\sqrt{E[D_{11}^2]} = O_P(\sqrt{s/\underline{C}})O(1)$. This along with the conclusion of Step 3 show $|\widehat{Q} - Q| = o_P(1)$. Applying the continuous mapping theorem under Assumption 2(1) then gives $|\widehat{Q}^{-1} - Q^{-1}| = o_P(1)$.

Step 5. Finally, $|\widehat{\Gamma}| \leq |\Gamma| + |\widehat{\Gamma} - \Gamma| = O_P(1)$ following the bounds from Steps 1 and 2 and Assumption 2 (1). Similarly, $|\widehat{Q}^{-1}| \leq |Q^{-1}| + |\widehat{Q}^{-1} - Q^{-1}|$ are both bounded following Assumption 2 (1) and Steps 3 and 4. \blacksquare

APPENDIX B. AUXILIARY LEMMAS

The following Lemma is an immediate consequence of Theorem 5.1 of Chernozhukov, Chetverikov and Kato (2014) and Lemma 8 of Chernozhukov, Chetverikov and Kato (2015).

Lemma 2 (A Concentration Inequality). *Let $(X_i)_{i \in [n]}$ be p -dimensional independent random vectors and let $B = \sqrt{E[\max_{i \in [n]} \|X_i\|_\infty^2]}$ and $\sigma^2 = \max_{j \in [p]} \frac{1}{n} \sum_{i=1}^n E|X_{ij}|^2$. Then with probability at least $1 - C(\log n)^{-1}$,*

$$\max_{j \in [p]} \left| \frac{1}{n} \sum_{i=1}^n (|X_{ij}| - E|X_{ij}|) \right| \lesssim \sqrt{\frac{\sigma^2 \log(p \vee n)}{n}} + \frac{B \log(p \vee n)}{n}.$$

The following is an immediate result of Lemma 10 of Belloni, Chen, Chernozhukov and Hansen (2012) with $n = NM$, $\lambda = CNM\sqrt{\log a/\underline{C}}$ for $C > 1$ and $c_s = \sqrt{s/\underline{C}}$.

Lemma 3 (Sparsity Bound for Lasso). *Consider lasso estimator (3.3) and suppose Assumption 3 (1)–(2) and 4. Then suppose $\lambda_2/NM \geq c\|(NM)^{-1} \sum_{i=1}^N \sum_{j=1}^M X_{ij}v_{ij}\|_\infty$ w.p. $1 - o(1)$, then denote $\widehat{s} = \text{support}(\widehat{\gamma})$, we have $\widehat{s} \lesssim_P s$. Similar result holds for lasso estimator (3.4) as well.*

REFERENCES

- Anderson, Siwan. “Legal origins and female HIV.” *American Economic Review*, 108, no. 6 (2018): 1407-39.
- Andrews, Donald W.K. “Cross-section regression with common shocks.” *Econometrica*, 73, no. 5 (2005): 1551-1585.
- Belloni, Alexandre, Daniel Chen, Victor Chernozhukov, and Christian Hansen. “Sparse models and methods for optimal instruments with an application to eminent domain.” *Econometrica*, 80, no. 6 (2012): 2369-2429.
- Belloni, Alexandre, and Victor Chernozhukov. “Least squares after model selection in high-dimensional sparse models.” *Bernoulli*, 19, no. 2 (2013): 521-547.
- Belloni, Alexandre, Victor Chernozhukov, and Christian Hansen. “Inference for high-dimensional sparse econometric models.” arXiv preprint arXiv:1201.0220 (2011).
- Belloni, Alexandre, Victor Chernozhukov, and Christian Hansen. “Inference on treatment effects after selection among high-dimensional controls.” *Review of Economic Studies*, 81, no. 2 (2014): 608-650.
- Belloni, Alexandre, Victor Chernozhukov, Christian Hansen, and Damian Kozbur. “Inference in high-dimensional panel models with an application to gun control.” *Journal of Business and Economic Statistics*, 34, no. 4 (2016): 590-605.
- Belloni, Alexandre, Victor Chernozhukov, and Kengo Kato. “Uniform post-selection inference for least absolute deviation regression and other Z-estimation problems.” *Biometrika*, 102, no. 1 (2015): 77-94.
- Cameron, Colin A., Jonah B. Gelbach, and Douglas L. Miller. “Robust inference with multiway clustering.” *Journal of Business and Economic Statistics*, 29, no. 2 (2011): 238-249.
- Cameron, Colin A. and Douglas L. Miller. “A practitioner’s guide to cluster-robust inference.” *Journal of Human Resources* 50, no. 2 (2015): 317-372.
- Chernozhukov, Victor, Denis Chetverikov, and Kengo Kato. “Gaussian approximation of suprema of empirical processes.” *Annals of Statistics*, 42, no. 4 (2014): 1564-1597.
- Chernozhukov, Victor, Denis Chetverikov, and Kengo Kato. “Comparison and anti-concentration bounds for maxima of Gaussian random vectors.” *Probability Theory and Related Fields* 162, no. 1-2 (2015): 47-70.

- Davezies, Laurent, Xavier D'Haultfoeuille, and Yannick Guyonvarch. "Asymptotic results under multiway clustering." arXiv preprint arXiv:1807.07925 (2018).
- Davezies, Laurent, Xavier D'Haultfoeuille, and Yannick Guyonvarch. "Empirical Process Results for Exchangeable Arrays." arXiv preprint arXiv:1906.11293 (2019).
- Dickens, Andrew. "Ethnolinguistic favoritism in African politics." *American Economic Journal: Applied Economics*, 10, no. 3 (2018): 370-402.
- Gershman, Boris. "Witchcraft beliefs and the erosion of social capital: Evidence from Sub-Saharan Africa and beyond." *Journal of Development Economics*, 120 (2016): 182-208.
- Javanmard, Adel and Andrea Montanari. "Confidence intervals and hypothesis testing for high-dimensional regression." *Journal of Machine Learning Research*, 15, no. 1 (2014): 2869-2909.
- Kallenberg, Olav. Probabilistic symmetries and invariance principles. Springer Science and Business Media (2005).
- Kock, Anders Bredahl. "Oracle inequalities, variable selection and uniform inference in high-dimensional correlated random effects panel data models." *Journal of Econometrics*, 195, no. 1 (2016): 71-85.
- MacKinnon, James G., Morten Ørregaard Nielsen, and Matthew D. Webb. "Wild Bootstrap and Asymptotic Inference with Multiway Clustering." No. 1415. Queen's Economics Department Working Paper, 2019.
- Menzel, Konrad. "Inference for games with many players." *Review of Economic Studies*, 83, no. 1 (2015): 306-337.
- Menzel, Konrad. "Bootstrap with clustering in two or more dimensions." arXiv preprint arXiv:1703.03043 (2017).
- Michalopoulos, Stelios and Elias Papaioannou. "PreColonial Ethnic Institutions and Contemporary African Development." *Econometrica*, 81, no. 1 (2013) 113-152.
- Michalopoulos, Stelios, and Elias Papaioannou. "National institutions and subnational development in Africa." *The Quarterly Journal of Economics*, 129, no. 1 (2013): 151-213.
- Michalopoulos, Stelios, and Elias Papaioannou. "The long-run effects of the scramble for Africa." *American Economic Review*, 106, no. 7 (2016): 1802-48.
- Nunn, Nathan and Leonard Wantchekon. "The slave trade and the origins of mistrust in Africa." *American Economic Review*, 101, no. 7 (2011): 3221-3252.

van de Geer, Sara, Peter Bühlmann, Ya'acov Ritov, and Ruben Dezeure. "On asymptotically optimal confidence regions and tests for high-dimensional models." *Annals of Statistics*, 42, no. 3 (2014): 1166-1202.

Zhang, Cun-Hui, and Stephanie S. Zhang. "Confidence intervals for low dimensional parameters in high dimensional linear models." *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76, no. 1 (2014): 217-242.

(Harold Chiang) DEPARTMENT OF ECONOMICS, VANDERBILT UNIVERSITY, UNITED STATES

(Yuya Sasaki) DEPARTMENT OF ECONOMICS, VANDERBILT UNIVERSITY, UNITED STATES

N	M	Dim	Statistics				95% Coverage		
			Avg	Bias	SD	RMSE	0-Way	1-Way	2-Way
10	10	100	0.491	-0.009	0.172	0.172	0.808	0.797	0.919
20	20	100	0.499	-0.001	0.076	0.076	0.855	0.858	0.964
40	40	100	0.501	0.001	0.045	0.045	0.792	0.848	0.959
10	10	200	0.480	-0.020	0.356	0.357	0.753	0.747	0.877
20	20	200	0.498	-0.002	0.075	0.075	0.863	0.860	0.962
40	40	200	0.500	0.000	0.041	0.041	0.830	0.859	0.962
10	10	400	0.459	-0.041	0.744	0.745	0.682	0.690	0.822
20	20	400	0.496	-0.004	0.079	0.079	0.846	0.841	0.951
40	40	400	0.500	0.000	0.037	0.037	0.869	0.876	0.970
10	10	800	0.438	-0.062	0.959	0.961	0.615	0.634	0.764
20	20	800	0.492	-0.008	0.086	0.086	0.808	0.803	0.930
40	40	800	0.499	-0.001	0.037	0.037	0.870	0.871	0.968
10	10	1600	0.394	-0.106	1.140	1.140	0.555	0.585	0.704
20	20	1600	0.487	-0.013	0.098	0.099	0.763	0.762	0.903
40	40	1600	0.498	-0.002	0.038	0.038	0.862	0.859	0.964

TABLE 1. Simulation results. The first three columns indicate the two-way sample sizes (N, M) and the dimension (Dim) of $(\alpha, \beta)'$. The next four columns report simulation statistics for $\tilde{\alpha}$, including the average (Avg), bias (Bias), standard deviation (SD), and root mean square error (RMSE). The last three columns report 95% coverage frequencies of α with the heteroskedasticity robust variance estimator (0-Way), the one-way cluster-robust variance estimator (1-Way), and our multi-way cluster-robust variance estimator (2-Way). The data generating parameters are set to $(\omega_1^x, \omega_2^x) = (0.25, 0.25)$, $(\omega_1^\varepsilon, \omega_2^\varepsilon) = (0.25, 0.25)$, and $\rho = 0.50$. The results are based on 25,000 Monte Carlo iterations for each row.

Variables		Number of Observations	Cluster Size		Original Estimates	Lasso Estimates
Y	D		N	M		
Trust of Neighbors	Slave Exports	20,027	185	1,257	-0.00068 (0.00015)	-0.00083 (0.00022)
Trust of Neighbors	Exports/ Area	20,027	185	1,257	-0.019 (0.005)	-0.025 (0.007)
Trust of Neighbors	Exports/ Population	17,644	157	1,214	-0.531 (0.147)	-0.684 (0.232)
Trust of Neighbors	Log Slave Exports	20,027	185	1,257	-0.037 (0.014)	-0.045 (0.021)
Trust of Neighbors	Log Exports/ Area	20,027	185	1,257	-0.159 (0.034)	-0.210 (0.050)
Trust of Neighbors	Log Exports/ Population	17,644	157	1,214	-0.743 (0.187)	-0.957 (0.304)

TABLE 2. Estimates of the effects of slave trade on mistrust in Africa. The first two columns indicate which measures of the dependent and explanatory variables are used. The next three columns show the number of observations, the number of ethnic groups (N), and the number of districts (M). The last two columns show the original estimates obtained under the prototype model by Nunn and Wantchekon (2011, Table 1) and corresponding lasso estimates obtained under more flexible model specification by our method.

Variables		No.	Cluster Size		Population	Original	Lasso
<i>Y</i>	<i>D</i>	Obs.	<i>N</i>	<i>M</i>	Density	Estimates	Estimates
Light	Jurisdictional	682	93	48	No	0.2794	0.2266
Density	Hierarchy					(0.0852)	(0.0797)
Light	Jurisdictional	682	93	48	Yes	0.1766	0.1649
Density	Hierarchy					(0.0501)	(0.0541)
Light	Political	682	93	48	No	0.5049	0.4158
Density	Centralization					(0.1573)	(0.1451)
Light	Political	682	93	48	Yes	0.3086	0.2985
Density	Centralization					(0.0972)	(0.1080)

TABLE 3. Estimates of the effects of pre-colonial institutions on regional development in Africa. The first two columns indicate which measures of the dependent and explanatory variables are used. The next three columns show the number of observations, the number of ethnic groups (N), and the number of districts (M). The next column indicates a control for population density. The last two columns show the original estimates obtained under the prototype model by Michalopoulos and Papaioannou (2013, Table 3) and corresponding lasso estimates obtained under more flexible model specification by our method.