

DISTRIBUTED MEAN DIMENSION REDUCTION THROUGH SEMI-PARAMETRIC APPROACHES

Zhengtian Zhu, Wangli Xu and Liping Zhu*

Renmin University of China

Abstract: In this study we recast the semiparametric mean dimension reduction approaches under a least squares framework. This changes the problem of recovering the central mean subspace into a series of problems of estimating slopes in linear regressions, and enables us to incorporate penalties to produce sparse solutions. We further adapt the semiparametric mean dimension reduction approaches to distributed settings in which massive data are scattered at various locations, and cannot be aggregated or processed by a single machine. We propose three communication efficient distributed algorithms. The first yields a dense solution, the second produces a sparse estimation, and the third provides an orthonormal basis. The distributed algorithms are less complex computationally than a pooled algorithm, and attain oracle rates after a finite number of iterations. Using extensive numerical studies, we demonstrate the finite sample performance of the distributed estimates, and compare it with that of a pooled algorithm.

Key words and phrases: Central subspace, distributed estimation, sufficient dimension reduction.

1. Introduction

Recent advances in science and technology allow us to collect and process massive data in a cost-efficient manner. However, while such data present significant opportunities, they also present challenges to statisticians and data scientists.

Massive data usually have a high dimension and a large volume. To cope with the high dimensionality, sufficient dimension reduction (Li, 1991; Cook, 2009) is an effective paradigm that combines the idea of a linear reduction with the notion of sufficiency. Cook and Li (2002) introduced the concept of mean dimension reduction, which concerns $E(Y | \mathbf{x})$, where Y is a univariate response and $\mathbf{x} = (X_1, \dots, X_p)^\top$ is a p -vector of covariates. The central mean subspace model assumes that there exists a $p \times d$ matrix $\boldsymbol{\beta} \in \mathbb{R}^{p \times d}$ such that

$$E(Y | \mathbf{x}) = E(Y | \mathbf{x}^\top \boldsymbol{\beta}). \quad (1.1)$$

The column space of $\boldsymbol{\beta}$, denoted by $\mathcal{S}(\boldsymbol{\beta})$, is referred to as the mean dimension reduction subspace. If the intersection of all such subspaces satisfies (1.1), we

*Corresponding author. E-mail: zhu.liping@ruc.edu.cn

call it the central mean subspace, which is unique and denoted as $\mathcal{S}_{E(Y|\mathbf{x})}$. The column dimension of $\boldsymbol{\beta}$, denoted by d , is an integer between zero and p . In the trivial case of $d = 0$, Y is mean independent of \mathbf{x} ; that is, $E(Y | \mathbf{x}) = E(Y)$. In the special case of $d = p$, we can simply set $\boldsymbol{\beta}$ to be the $p \times p$ identity matrix $\mathbf{I}_{p \times p}$, and model (1.1) is always true. In general, d must be decided using a data-driven mechanism.

Many approaches have been proposed to identify and recover $\mathcal{S}_{E(Y|\mathbf{x})}$. These approaches can be roughly classified into three categories. The first consists mainly of inverse regression methods, and requires stringent distributional assumptions on the covariates, such as the linearity mean (Li, 1991) and constant variance conditions (Cook and Weisberg, 1991). Examples include the ordinary least squares (Li and Duan, 1989), principal Hessian directions (Li, 1992), and their variations (Cook and Li, 2004). The second category consists mainly of forward regression methods, which extract the information of $\mathcal{S}_{E(Y|\mathbf{x})}$ from the derivatives of $E(Y | \mathbf{x})$. Examples include the average derivative estimation (Härdle and Stoker, 1989) and minimum average variance estimation (Xia et al., 2002). The third category comprises the semiparametric estimating equations approaches, which require minimal distributional assumptions on the covariates. Here, examples include the works of Ma and Zhu (2014), Luo, Li and Yin (2014), and Zhu and Zhong (2015). In particular, Luo, Li and Yin (2014) thoroughly examine the asymptotic properties of the semiparametric approaches of Ma and Zhu (2014) when the variance function $\text{var}(Y | \mathbf{x})$ is estimated consistently using kernel smoothers. Zhu and Zhong (2015) extend the work of Ma and Zhu (2014) by allowing for multiple responses, and assume implicitly that the variance functions are all constants.

We advocate using the semiparametric approaches of Ma and Zhu (2014), for at least two reasons. First, they avoid having to use the linearity mean and constant variance conditions, thus generalizing the usefulness of sufficient dimension reduction. Indeed, even when these distributional assumptions are satisfied, sufficient dimension reduction methods with the linearity mean and constant variance estimated using nonparametric treatments are more efficient than those that use these conditions directly (Ma and Zhu, 2013). Second, the semiparametric approaches are locally efficient. Specifically, the resultant solutions attain the semiparametric efficiency bound as long as the variance function $\text{var}(Y | \mathbf{x})$ is specified correctly (Ma and Zhu, 2014) or estimated consistently (Luo, Li and Yin, 2014), and remain consistent even when $\text{var}(Y | \mathbf{x})$ is misspecified or estimated inconsistently.

The problem with the semiparametric approaches is that they have computational complexity of $O(N^2)$, where N is the total sample size, and thus are usually computationally prohibitive for massive data. Massive data are often scattered across various locations, possibly because of memory or storage limitations, or privacy concerns. Therefore, to cope with large volumes of massive

data, distributed methodologies are highly desirable.

Although distributed statistical inference has received considerable attention, most existing approaches require only one round of communication: the node machines conduct the inference in parallel and send the results to the central machine, which aggregates all information to produce a final solution; see, for example, Zhang, Duchi and Wainwright (2013), Battey et al. (2018), and Fan et al. (2019). While these one-shot methods are communication efficient, they only work with a small number of node machines, and require large sample size on each of them. Violating these requirements results in suboptimal performance. Balcan et al. (2016) designed a distributed algorithm for a kernel principal component analysis. They obtain approximate solutions with a relatively low communication cost. Jordan, Lee and Yang (2019) and Fan, Guo and Wang (2023) developed iterative methods with multiple rounds of aggregations, which substantially relaxes the requirement on the number of machines. Cai, Li and Zhu (2020) proposed implementing a sliced inverse regression in an online manner, if the observations arrive in data streams. Chen, Xu and Zhu (2022) and Zhu and Zhu (2022) adapt sufficient dimension reduction with convex loss functions to distributed settings. The resultant estimates possess nearly oracle rates after a finite number of iterations. However, both require stringent distributional assumptions, such as the linearity mean condition. In addition, few works have examined how to conduct distributed statistical inference in the context of semiparametric estimating equations.

In this study, we reformulate the Newton–Raphson iterations of the semiparametric estimating equations under a least squares framework. This changes the problem of estimating the central mean subspace into a series of problems of estimating slopes in linear regressions, enabling us to incorporate penalties to yield sparse solutions. We propose three distributed algorithms, under various identifiability conditions. The first algorithm yields a dense solution, the second produces a sparse estimation, and the third provides an orthonormal basis.

Our proposed distributed algorithms possess at least three desirable properties. First, these algorithms are communication efficient. The estimating equations themselves and their gradients correspond to the gradients and Hessians, respectively, of the least squares losses. We estimate the gradients separately using the observations recorded in each node machine. The gradients are then transmitted to the central node to form an aggregated estimation of the overall gradient. The communication cost is $O(mp)$, where m is the number of node machines. This is the minimal price we have to pay in distributed settings. Instead of using all N observations scattered across m node machines, we estimate the Hessians simply using the observations from the central node machine only, which incurs no transmission cost. In this sense, our proposed distributed algorithms are communication efficient. Second, the resultant distributed estimates possess desirable theoretical properties. For

example, they achieve the oracle rate after a finite number of iterations. We derive the contraction rate for the distributed estimates. After a small number of iterations, the optimization errors are asymptotically negligible compared with the statistical errors. Therefore, in an asymptotic sense, the distributed estimates behave as well as the classic pooled estimate, which requires pooling all observations on in a single machine. Lastly, the distributed algorithms are computationally much more efficient than the corresponding pooled algorithm.

The remainder of the paper is organized as follows. In Section 2, we review the semiparametric approaches of Ma and Zhu (2014), and recast them into a least squares framework. Note that, although $\mathcal{S}_{E(Y|\mathbf{x})}$ is unique, its basis matrix β is not. Two sets of conditions ensure that β is identifiable. One requires that the upper $d \times d$ block of β is the identity matrix $\mathbf{I}_{d \times d}$, and thus the lower $(p - d) \times d$ block comprises free parameters. The second requires β to be orthonormal, that is, $\beta^T \beta = \mathbf{I}_{d \times d}$. This condition is widely used, although it is not sufficient to ensure the identifiability of β unless some additional assumptions are imposed. In Sections 3 and 4, we adapt the semiparametric approaches to distributed settings under the first set of identifiability conditions. In Section 5, we suggest a distributed algorithm under the orthogonality constraint. We discuss our simulation studies in Section 6, and conclude our paper in Section 7.

2. A Brief Review of Semiparametric Approaches and Equivalent Reformulations

2.1. Notations

Let $C, C_0, C_1, \dots, c, c_0, c_1, \dots$ be generic constants that may vary at each appearance. For a vector $\alpha = (\alpha_1, \dots, \alpha_p)^T$, we define $|\alpha|_1 \stackrel{\text{def}}{=} \sum_{i=1}^p |\alpha_i|$ and $|\alpha|_2 \stackrel{\text{def}}{=} (\sum_{i=1}^p \alpha_i^2)^{1/2}$. For a matrix $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{p \times d}$, $|\mathbf{A}|_\infty \stackrel{\text{def}}{=} \max_{1 \leq i \leq p, 1 \leq j \leq d} |a_{ij}|$ and $\|\mathbf{A}\|_\infty \stackrel{\text{def}}{=} \max_{1 \leq i \leq p} \sum_{1 \leq j \leq d} |a_{ij}|$. Furthermore, $\text{vec}(\mathbf{A})$ is an operator that stacks all columns of \mathbf{A} vertically in order, and $\text{vecl}(\mathbf{A})$ is an operator that vectorizes the lower $(p - d) \times d$ block of \mathbf{A} , that is, $\text{vecl}(\mathbf{A}) = \text{vec}(\mathbf{A}_2)$, for $\mathbf{A} = (\mathbf{A}_1^T, \mathbf{A}_2^T)^T$, $\mathbf{A}_1 \in \mathbb{R}^{d \times d}$, and $\mathbf{A}_2 \in \mathbb{R}^{(p-d) \times d}$. In addition, $\text{vec}^{-1}(\cdot)$ is an inverse operator of $\text{vec}(\cdot)$, that rearranges a (pd) -vector as a $p \times d$ matrix in column order. In particular, $\text{vec}^{-1}\{\text{vec}(\mathbf{A})\} = \mathbf{A}$. The largest and smallest singular values of \mathbf{A} are denoted by $\lambda_{\max}(\mathbf{A})$ and $\lambda_{\min}(\mathbf{A})$, respectively. Let $\mathbf{P}(\mathbf{A}) \stackrel{\text{def}}{=} \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ be the projection matrix of \mathbf{A} . Define $\mathbf{A}^{\otimes 2} = \mathbf{A} \mathbf{A}^T$, and let $\mathbf{A} \otimes \mathbf{B}$ be the Kronecker product of \mathbf{A} and \mathbf{B} , where $\mathbf{B} \in \mathbb{R}^{p_1 \times d_1}$. For two sequences of real numbers, $\{a_n\}_{n=1}^\infty$ and $\{b_n\}_{n=1}^\infty$, we write $a_n = O(b_n)$ if there exists a positive constant C such that $|a_n/b_n| \leq C$, for sufficiently large n . For two sequences of random variables, $\{X_n\}_{n=1}^\infty$ and $\{Y_n\}_{n=1}^\infty$, we write $X_n = O_p(Y_n)$, if, for any $\varepsilon > 0$, there exists $C > 0$ such that $\text{pr}(|X_n/Y_n| \leq C) \geq 1 - \varepsilon$, for sufficiently large n .

2.2. Recasting semiparametric approaches under a least squares framework

Here, we briefly review the semiparametric approaches of Ma and Zhu (2014). With a slight abuse of notation, we denote as β the basis matrix of $\mathcal{S}_{E(Y|\mathbf{x})}$, with its upper $d \times d$ block being $\mathbf{I}_{d \times d}$, and all other elements of β , $\text{vecl}(\beta)$, being free parameters. Let $m(\mathbf{x}^T \beta) \stackrel{\text{def}}{=} E(Y \mid \mathbf{x}^T \beta)$, $\mathbf{m}_1(\mathbf{x}^T \beta) \stackrel{\text{def}}{=} \text{vec}\{\partial m(\mathbf{x}^T \beta) / \partial (\mathbf{x}^T \beta)\}$, $\varepsilon \stackrel{\text{def}}{=} Y - m(\mathbf{x}^T \beta)$, and $w(\mathbf{x}) \stackrel{\text{def}}{=} \{E(\varepsilon^2 \mid \mathbf{x})\}^{-1}$. Let $\alpha \in \mathbb{R}^{p \times d}$ be an intermediate estimate, with its upper $d \times d$ block being $\mathbf{I}_{d \times d}$. Let $\text{vecl}(\alpha) \in \mathbb{R}^{(p-d)d \times 1}$ be a vector of free parameters. Define

$$\tilde{\mathbf{x}}(\alpha) \stackrel{\text{def}}{=} \text{vecl} \left\{ \left[\mathbf{x} - \frac{E\{\mathbf{x}w(\mathbf{x}) \mid \mathbf{x}^T \alpha\}}{E\{w(\mathbf{x}) \mid \mathbf{x}^T \alpha\}} \right] \mathbf{m}_1^T(\mathbf{x}^T \alpha) \right\} \in \mathbb{R}^{(p-d)d \times 1}. \quad (2.1)$$

We further write $\mathbf{S}\{\mathbf{x}, Y, \alpha, w(\mathbf{x})\} \stackrel{\text{def}}{=} \{Y - m(\mathbf{x}^T \alpha)\}w(\mathbf{x})\tilde{\mathbf{x}}(\alpha)$. Ma and Zhu (2014) showed that $E[\mathbf{S}\{\mathbf{x}, Y, \beta, w(\mathbf{x})\}] = \mathbf{0}$. In other words, solving the estimating equations yields a consistent estimate for the basis β of $\mathcal{S}_{E(Y|\mathbf{x})}$.

We seek β by using Newton–Raphson iterations. This requires calculating the gradient of $E[\mathbf{S}\{\mathbf{x}, Y, \alpha, w(\mathbf{x})\}]$ with respect to $\text{vecl}(\alpha)$, which yields $\{-\mathbf{H}(\alpha)\}$, where $\mathbf{H}(\alpha) \stackrel{\text{def}}{=} E[w(\mathbf{x})\{\tilde{\mathbf{x}}(\alpha)\}\{\tilde{\mathbf{x}}(\alpha)\}^T]$. Start from an initial value $\beta^{(0)}$. The Newton–Raphson iteration proceeds as

$$\text{vecl}(\beta^{(t+1)}) \stackrel{\text{def}}{=} \text{vecl}(\beta^{(t)}) + \left\{ \mathbf{H}(\beta^{(t)}) \right\}^{-1} E\left[\mathbf{S}\{\mathbf{x}, Y, \beta^{(t)}, w(\mathbf{x})\}\right]. \quad (2.2)$$

Throughout, we fix the upper $d \times d$ block of $\beta^{(t)}$ to be $\mathbf{I}_{d \times d}$. We update $\beta^{(t)}$ with $\beta^{(t+1)}$, and iterate (2.2) until convergence.

Next, we present our first contribution to the literature, where we recast the above Newton–Raphson iteration under a least squares framework. Define

$$\tilde{Y}(\alpha) \stackrel{\text{def}}{=} \{\tilde{\mathbf{x}}(\alpha)\}^T \text{vecl}(\alpha) + \{Y - m(\mathbf{x}^T \alpha)\}. \quad (2.3)$$

Here, (2.2) can be written equivalently as

$$\text{vecl}(\beta^{(t+1)}) = \left\{ \mathbf{H}(\beta^{(t)}) \right\}^{-1} E \left\{ w(\mathbf{x}) \tilde{\mathbf{x}}(\beta^{(t)}) \tilde{Y}(\beta^{(t)}) \right\},$$

which exactly minimizes the following weighted least squares loss function:

$$\text{vecl}(\beta^{(t+1)}) = \underset{\alpha}{\text{argmin}} E \left[\{\tilde{Y}(\beta^{(t)}) - \tilde{\mathbf{x}}(\beta^{(t)})^T \text{vecl}(\alpha)\}^2 w(\mathbf{x}) \right], \quad (2.4)$$

for $t \geq 0$. We update $\beta^{(t)}$ with $\beta^{(t+1)}$, iterate (2.4) to obtain $\beta^{(t+2)}$, and so on. This iteration proceeds until convergence.

In the following, we show that this reformulation enables us to incorporate penalties on the right-hand side of (2.4), which are introduced to produce

sparse solutions. This equivalent reformulation is thus very appealing in high dimensions. In what follows, we introduce three distributed algorithms. The first yields a dense solution, the second produces a sparse solution, and the third provides an orthonormal estimate. These distributed algorithms are communication efficient, and the resultant solutions possess desirable theoretical properties. To save space, we relegate the description of the pooled algorithm to the online Supplementary Material.

3. The First Distributed Algorithm with Dense Solutions

3.1. The first communication efficient distributed algorithm

First, we explore how to adapt the above Newton–Raphson iterations to distributed settings when the observations are scattered across various locations. With a slight abuse of notation, we denote the observations $\{(\mathbf{x}_i, Y_i), i = 1, \dots, N\}$ as $\{(\mathbf{x}_{i,j}, Y_{i,j}), i = 1, \dots, n, j = 1, \dots, m\}$, assuming they are scattered across m machines. We further assume that the total sample size $N = nm$ is so large that a single machine cannot process all observations simultaneously, owing to memory or storage limitations. In this case, we require a communication efficient distributed algorithm.

Instead of using all N observations to estimate $m(\mathbf{x}^T \boldsymbol{\alpha})$, $\mathbf{m}_1(\mathbf{x}^T \boldsymbol{\alpha})$, $E\{w(\mathbf{x}) \mid \mathbf{x}^T \boldsymbol{\alpha}\}$, and $E\{\mathbf{x}w(\mathbf{x}) \mid \mathbf{x}^T \boldsymbol{\alpha}\}$, we suggest estimating them using the observations in the j th machine only, which yields m distinct estimates. In particular, we define $(\hat{b}_{k,j}, \hat{\mathbf{b}}_{k,j}) \stackrel{\text{def}}{=} \arg\min_{b_{k,j}, \mathbf{b}_{k,j}} \sum_{i=1, i \neq k}^n \{Y_{i,j} - b_{k,j} - (\mathbf{x}_{i,j}^T \boldsymbol{\alpha} - \mathbf{x}_{k,j}^T \boldsymbol{\alpha})\mathbf{b}_{k,j}\}^2 K_{h_1}(\mathbf{x}_{i,j}^T \boldsymbol{\alpha} - \mathbf{x}_{k,j}^T \boldsymbol{\alpha})$.

$$\arg\min_{b_{k,j}, \mathbf{b}_{k,j}} \sum_{i=1, i \neq k}^n \{Y_{i,j} - b_{k,j} - (\mathbf{x}_{i,j}^T \boldsymbol{\alpha} - \mathbf{x}_{k,j}^T \boldsymbol{\alpha})\mathbf{b}_{k,j}\}^2 K_{h_1}(\mathbf{x}_{i,j}^T \boldsymbol{\alpha} - \mathbf{x}_{k,j}^T \boldsymbol{\alpha}).$$

Let $\hat{m}_j(\mathbf{x}_{k,j}^T \boldsymbol{\alpha}) = \hat{b}_{k,j}$ and $\hat{\mathbf{m}}_{1,j}(\mathbf{x}_{k,j}^T \boldsymbol{\alpha}) = \hat{\mathbf{b}}_{k,j}$. In addition, we define

$$\hat{E}_j\{w(\mathbf{x}_{k,j}) \mid \mathbf{x}_{k,j}^T \boldsymbol{\alpha}\} \stackrel{\text{def}}{=} \frac{\sum_{i=1, i \neq k}^n K_{h_2}(\mathbf{x}_{i,j}^T \boldsymbol{\alpha} - \mathbf{x}_{k,j}^T \boldsymbol{\alpha})w(\mathbf{x}_{i,j})}{\sum_{i=1, i \neq k}^n K_{h_2}(\mathbf{x}_{i,j}^T \boldsymbol{\alpha} - \mathbf{x}_{k,j}^T \boldsymbol{\alpha})},$$

and

$$\hat{E}_j\{\mathbf{x}_{k,j}w(\mathbf{x}_{k,j}) \mid \mathbf{x}_{k,j}^T \boldsymbol{\alpha}\} \stackrel{\text{def}}{=} \frac{\sum_{i=1, i \neq k}^n K_{h_3}(\mathbf{x}_{i,j}^T \boldsymbol{\alpha} - \mathbf{x}_{k,j}^T \boldsymbol{\alpha})\{\mathbf{x}_{i,j}w(\mathbf{x}_{i,j})\}}{\sum_{i=1, i \neq k}^n K_{h_3}(\mathbf{x}_{i,j}^T \boldsymbol{\alpha} - \mathbf{x}_{k,j}^T \boldsymbol{\alpha})}.$$

Accordingly, we define

$$\hat{\mathbf{x}}_{k,j}(\boldsymbol{\alpha}) \stackrel{\text{def}}{=} \text{vecl} \left\{ \left[\mathbf{x}_{k,j} - \frac{\hat{E}_j\{\mathbf{x}_{k,j}w(\mathbf{x}_{k,j}) \mid \mathbf{x}_{k,j}^T \boldsymbol{\alpha}\}}{\hat{E}_j\{\mathbf{x}_{k,j} \mid \mathbf{x}_{k,j}^T \boldsymbol{\alpha}\}} \right] \hat{\mathbf{m}}_1^T(\mathbf{x}_{k,j}^T \boldsymbol{\alpha}) \right\}. \quad (3.1)$$

Define $\widehat{\mathbf{S}}_j\{\mathbf{x}_{k,j}, Y_{k,j}, \boldsymbol{\alpha}, w(\mathbf{x}_{k,j})\} \stackrel{\text{def}}{=} \{Y_{k,j} - \widehat{m}_j(\mathbf{x}_{k,j}^T \boldsymbol{\alpha})\} w(\mathbf{x}_{k,j}) \widehat{\mathbf{x}}_{k,j}(\boldsymbol{\alpha})$, and

$$\widehat{E}_j\left[\mathbf{S}\{\mathbf{x}_j, Y_j, \boldsymbol{\alpha}, w(\mathbf{x}_j)\}\right] \stackrel{\text{def}}{=} n^{-1} \sum_{k=1}^n \widehat{\mathbf{S}}_j\{\mathbf{x}_{k,j}, Y_{k,j}, \boldsymbol{\alpha}, w(\mathbf{x}_{k,j})\},$$

for $j = 1, \dots, m$. All are consistent estimates of $E[\mathbf{S}\{\mathbf{x}, Y, \boldsymbol{\alpha}, w(\mathbf{x})\}]$. Each has computational complexity of $O(n^2)$. More importantly, computing these quantities can be parallelized to further improve the computational efficiency. Therefore, as long as n is small relative to N , the computational complexity is reduced to $O(n^2)$ from $O(N^2)$, which is substantial. The above estimate, $\widehat{E}_j[\mathbf{S}\{\mathbf{x}_j, Y_j, \boldsymbol{\alpha}, w(\mathbf{x}_j)\}]$, uses only the observations from the j th machine, and thus can be computed in parallel. We transmit these estimates to the first central machine to form

$$\widehat{E}_{\text{dist},1}\left[\mathbf{S}\{\mathbf{x}, Y, \boldsymbol{\alpha}, w(\mathbf{x})\}\right] \stackrel{\text{def}}{=} m^{-1} \sum_{j=1}^m \widehat{E}_j\left[\mathbf{S}\{\mathbf{x}_j, Y_j, \boldsymbol{\alpha}, w(\mathbf{x}_j)\}\right],$$

which serves as an estimate of $E[\mathbf{S}\{\mathbf{x}, Y, \boldsymbol{\alpha}, w(\mathbf{x})\}]$. The communication cost of transmitting these quantities is $O(mp)$, which is the minimal price we have to pay in a distributed setting.

To implement (2.2), it remains to estimate $\mathbf{H}(\boldsymbol{\alpha})$. We use $\{(\mathbf{x}_{i,1}, Y_{i,1}), i = 1, \dots, n\}$, that is the observations from the first machine only. Specifically,

$$\widehat{\mathbf{H}}_j(\boldsymbol{\alpha}) \stackrel{\text{def}}{=} n^{-1} \sum_{k=1}^n w(\mathbf{x}_{k,j}) \widehat{\mathbf{x}}_{k,j}(\boldsymbol{\alpha}) \widehat{\mathbf{x}}_{k,j}^T(\boldsymbol{\alpha}), \text{ for } j = 1, \dots, m.$$

We implement the Newton–Raphson algorithm on the first machine. As such, there is no communication cost when estimating $\mathbf{H}(\boldsymbol{\alpha})$.

We propose the first communication efficient algorithm, which yields a dense solution. This is our second contribution to the literature. We start from an initial value $\boldsymbol{\beta}_{\text{dist},1}^{(0)}$, and then iterate the Newton–Raphson algorithm in a distributed fashion, as follows: $\text{vecl}(\boldsymbol{\beta}_{\text{dist},1}^{(t+1)}) \stackrel{\text{def}}{=}$

$$\text{vecl}(\boldsymbol{\beta}_{\text{dist},1}^{(t)}) + \left\{ \widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},1}^{(t)}) \right\}^{-1} \widehat{E}_{\text{dist},1}\left[\mathbf{S}\{\mathbf{x}, Y, \boldsymbol{\beta}_{\text{dist},1}^{(t)}, w(\mathbf{x})\}\right]. \quad (3.2)$$

Once we have $\boldsymbol{\beta}_{\text{dist},1}^{(t+1)}$ from the first machine, we update $\widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},1}^{(t)})$ with $\widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},1}^{(t+1)})$. Next, we broadcast $\boldsymbol{\beta}_{\text{dist},1}^{(t+1)}$ from the first machine to the remaining $(m-1)$ machines to update $\widehat{E}_j[\mathbf{S}\{\mathbf{x}_j, Y_j, \boldsymbol{\beta}_{\text{dist},1}^{(t)}, w(\mathbf{x}_j)\}]$ with $\widehat{E}_j[\mathbf{S}\{\mathbf{x}_j, Y_j, \boldsymbol{\beta}_{\text{dist},1}^{(t+1)}, w(\mathbf{x}_j)\}]$. The latter is transmitted to the first machine to form $\widehat{E}_{\text{dist},1}[\mathbf{S}\{\mathbf{x}, Y, \boldsymbol{\beta}_{\text{dist},1}^{(t+1)}, w(\mathbf{x})\}]$. We iterate (3.2) until convergence, and denote the final solution as $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$.

In the above distributed algorithm, we assume $w(\mathbf{x}_j)$ is known, which is unrealistic in practice. However, this is not problematic, because we can specify

$w(\mathbf{x}_j)$ as $w^*(\mathbf{x}_j)$, or assume it has a parametric form $w_j(\mathbf{x}_j, \theta_j)$. We can also estimate $w(\mathbf{x}_j)$ using a kernel smoother at each local machine. Specifically, at the j th local node, we estimate $w(\mathbf{x}_j)$ as

$$\widehat{w}_j(\mathbf{x}_{k,j}) \stackrel{\text{def}}{=} \frac{\sum_{i=1}^n K_{h_4}(\mathbf{x}_{i,j} - \mathbf{x}_{k,j})}{\sum_{i=1}^n K_{h_4}(\mathbf{x}_{i,j} - \mathbf{x}_{k,j}) \{Y_{i,j} - \widehat{m}_j(\mathbf{x}_{i,j}^T \boldsymbol{\alpha})\}^2}.$$

The consistency of $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$ does not depend on how we specify or estimate $w(\mathbf{x})$. In the following, we show that as long as $w(\mathbf{x})$ is specified correctly or estimated consistently, the distributed estimate $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$ is semiparametrically efficient, despite the convergence rate of $\widehat{w}_j(\mathbf{x})$ being slow in high dimensions. Even if $w(\mathbf{x})$ is incorrectly specified or inconsistently estimated, $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$ still possesses an oracle rate. To distinguish these distributed estimates, we write $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$ as $\widehat{\boldsymbol{\beta}}_{\text{dist},1}(w^*)$ if $w^*(\mathbf{x})$ is used, and as $\widehat{\boldsymbol{\beta}}_{\text{dist},1}(\widehat{w})$ if $\widehat{w}_j(\mathbf{x})$ is used on each local machine. When $w^*(\mathbf{x})$ is equal to $w(\mathbf{x})$, we denote the resulting distributed estimate $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$ as $\widehat{\boldsymbol{\beta}}_{\text{dist},1}(w)$. When this does not cause any ambiguity, we simply use $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$.

3.2. Theoretical properties of the first distributed algorithm

In the above distributed algorithm, we estimate $E[\mathbf{S}\{\mathbf{x}, Y, \boldsymbol{\alpha}, w(\mathbf{x})\}]$ using a divide-and-conquer strategy, and estimate $\mathbf{H}(\boldsymbol{\alpha})$ using the observations from the first machine only. This distributed is computationally much more efficient than the pooled algorithm. It is thus natural to ask whether the distributed estimate, $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$, is as “good” as the pooled estimate, $\widehat{\boldsymbol{\beta}}_{\text{pool},1}$, which amounts to studying the theoretical properties of $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$. Ma and Zhu (2014), Luo, Li and Yin (2014), and Luo and Cai (2016) study the theoretical properties of $\widehat{\boldsymbol{\beta}}_{\text{pool},1}$ thoroughly. We first present regularity conditions to establish the theoretical properties of $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$. Throughout, we suppose that the covariates \mathbf{x} and the response Y are centered, that is $E(\mathbf{x}) = \mathbf{0}$ and $E(Y) = 0$. Suppose d is a fixed number. We introduce the following regularity conditions to establish the theoretical result for $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$:

(C1) The Kernels: The multivariate kernel is a multiplication of univariate and symmetric kernels. The q th-order univariate kernel $K(\cdot)$ satisfies

$$\int K(u)du = 1, \int u^i K(u)du = 0, \quad 1 \leq i \leq q-1, \quad 0 \neq \int u^q K(u)du < \infty.$$

It has a compact support over which it is Lipschitz continuous.

(C2) The Density: The density function of $(\mathbf{x}^T \boldsymbol{\beta})$, denoted by $f(\mathbf{x}^T \boldsymbol{\beta})$, and $w(\mathbf{x}) \stackrel{\text{def}}{=} \{E(\varepsilon^2 | \mathbf{x})\}^{-1}$ are bounded away from zero and infinity.

(C3) The Smoothness: Let $r(\mathbf{x}^T \boldsymbol{\alpha}) \stackrel{\text{def}}{=} E\{a(\mathbf{x}, Y)f(\mathbf{x}^T \boldsymbol{\alpha}) | \mathbf{x}^T \boldsymbol{\alpha}\}$, for $a(\mathbf{x}, Y)$ being Y , $w(\mathbf{x})$, or $\mathbf{x}w(\mathbf{x})$. The $(q-1)$ th derivatives of $r(\mathbf{x}^T \boldsymbol{\alpha})$, $f(\mathbf{x}^T \boldsymbol{\alpha})$, and $m(\mathbf{x}^T \boldsymbol{\alpha})$ are Lipschitz continuous in the neighborhood of $(\mathbf{x}^T \boldsymbol{\beta})$.

(C4) The Covariate: The covariate \mathbf{x} is sub-Gaussian. Let $\Sigma \stackrel{\text{def}}{=} \text{cov}(\mathbf{x}, \mathbf{x}^T)$. There exists $c > 1$ such that $c^{-1} \leq \lambda_{\min}(\Sigma) \leq \lambda_{\max}(\Sigma) \leq c$. $\mathbf{H}(\boldsymbol{\alpha})$ is invertible at $\boldsymbol{\beta}$. There exists $\lambda > 0$ such that $\mathbf{H}(\boldsymbol{\beta}) \geq \lambda \mathbf{I}_{(p-d)d \times (p-d)d}$.

(C5) The Moments: $E(\|\mathbf{x}\|_2^4) < \infty$, $E(Y^4) < \infty$, and $E\{\|\mathbf{m}_1(\mathbf{x}^T \boldsymbol{\beta})\|_2^4\} < \infty$. There exist G and H such that $E\{|\widehat{E}_j[\mathbf{S}\{\mathbf{x}, Y, \boldsymbol{\alpha}, w(\mathbf{x})\}]\|_2^4\} \leq G^4$ and $E\{\|\widehat{\mathbf{H}}(\boldsymbol{\alpha}) - \mathbf{H}(\boldsymbol{\alpha})\|_2^4\} \leq H^4$. There exists $L(\mathbf{x}, Y)$ such that $\|\widehat{\mathbf{H}}(\boldsymbol{\alpha}_1; \mathbf{x}) - \widehat{\mathbf{H}}(\boldsymbol{\alpha}_2; \mathbf{x})\|_2 \leq L(\mathbf{x}, Y)\|\boldsymbol{\alpha}_1 - \boldsymbol{\alpha}_2\|_2$, for $E\{L(\mathbf{x}, Y)\} \leq L^4$.

(C6) The Bandwidths: The bandwidths satisfy $Nh_k^{2q}h_l^{2q} \rightarrow 0$, $Nh_1^{2(q-1)}h_l^{2q} \rightarrow 0$, and $Nh_k^d h_l^d \rightarrow \infty$, for $1 \leq k \leq l \leq 4$.

(C7) The Sample Size: There exist $c_1 > 0$ and $0 < c_2 < 1$ such that $n \geq \max(c_1 p, N^{c_2})$.

(C8) The Initial Value: The initial value satisfies $|\text{vecl}(\boldsymbol{\beta}^{(0)} - \boldsymbol{\beta})|_2 = O_p(n^{-1/2})$.

We assume these conditions for technical reasons, though they are widely assumed in literature. In particular, condition (C1) allows for the second-order kernels, and condition (C6) allows for optimal bandwidths. The distributed algorithm requires the sample size n to be large enough to satisfy condition (C7), which is required to ensure condition (C8) holds when we calculate the initial value $\boldsymbol{\beta}^{(0)}$ using the observations from the first machine only. Condition (C7) is also used by Zhang, Duchi and Wainwright (2013), Battey et al. (2015), and Jordan, Lee and Yang (2019), and is typically regarded as mild in distributed settings. Condition (C8) requires that the initial value be consistent.

We provide a high probability error bound for $\boldsymbol{\beta}_{\text{dist},1}^{(t)}$ in the following theorem.

Theorem 1. *Under conditions (C1)–(C8), we have for $t \geq 1$,*

$$|\text{vecl}(\boldsymbol{\beta}_{\text{dist},1}^{(t)} - \boldsymbol{\beta})|_2 = O_p\left\{n^{-(t+1)/2} + N^{-1/2}\right\}.$$

An important implication of Theorem 1 is that $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$ can behave as well as $\widehat{\boldsymbol{\beta}}_{\text{pool},1}$ only after a finite number of iterations. In order to ensure $|\text{vecl}(\boldsymbol{\beta}_{\text{dist},1}^{(t+1)} - \widehat{\boldsymbol{\beta}}_{\text{pool},1})|_2 = o_p(N^{-1/2})$, we are merely required to conduct at most $\lceil \log N / \log n \rceil$ iterations, where $\lceil x \rceil$ denotes the smallest integer larger than or equal to x . In other words, for a sufficiently large t , the optimization error of the distributed estimate, $|\text{vecl}(\widehat{\boldsymbol{\beta}}_{\text{dist},1} - \widehat{\boldsymbol{\beta}}_{\text{pool},1})|_2$, is almost negligible in comparison with the statistical error of the pooled estimate, $|\text{vecl}(\widehat{\boldsymbol{\beta}}_{\text{pool},1} - \boldsymbol{\beta})|_2$. That is, $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$ behaves as well as $\widehat{\boldsymbol{\beta}}_{\text{pool},1}$. If N is a polynomial order of n , $\lceil \log N / \log n \rceil$ is a finite number. In other words, we are required to conduct at most a finite number of iterations to ensure that $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$ is almost as good as $\widehat{\boldsymbol{\beta}}_{\text{pool},1}$. As a consequence, $\widehat{\boldsymbol{\beta}}_{\text{dist},1}$ shares almost the same theoretical properties as $\widehat{\boldsymbol{\beta}}_{\text{pool},1}$.

4. The Second Distributed Algorithm with Sparse Solutions

4.1. The second communication efficient algorithm

Next, we introduce a distributed algorithm that yields a sparse solution. Under the least squares framework (2.4), we can incorporate penalties into the loss function to produce sparse solutions. In particular, we define

$$\widehat{Y}(\boldsymbol{\alpha}) \stackrel{\text{def}}{=} \{\widehat{\mathbf{x}}(\boldsymbol{\alpha})\}^T \text{vecl}(\boldsymbol{\alpha}) + \{Y - \widehat{m}(\mathbf{x}^T \boldsymbol{\alpha})\}. \quad (4.1)$$

We incorporate the least absolute shrinkage and selection operator (Tibshirani, 1996) into the least squares framework. With the observations denoted as $\{(\mathbf{x}_{i,k}, Y_{i,k}), i = 1, \dots, n, k = 1, \dots, m\}$, we ignore the penalty for now, and rewrite the least squares loss function as follows:

$$\mathcal{L}_N(\boldsymbol{\alpha}) \stackrel{\text{def}}{=} (2N)^{-1} \sum_{i=1}^n \sum_{k=1}^m \{\widehat{Y}_{i,k}(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) - \widehat{\mathbf{x}}_{i,k}(\boldsymbol{\beta}_{\text{dist},2}^{(t)})^T \text{vecl}(\boldsymbol{\alpha})\}^2 w(\mathbf{x}_{i,k}),$$

where $\boldsymbol{\beta}_{\text{dist},2}^{(t)}$ is an intermediate distributed estimate. We define

$$\begin{aligned} \mathbf{z}_k(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) &\stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n w(\mathbf{x}_{i,k}) \widehat{\mathbf{x}}_{i,k}(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) \widehat{Y}_{i,k}(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) \text{ and} \\ \mathbf{z}_N(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) &\stackrel{\text{def}}{=} \frac{1}{m} \sum_{k=1}^m \mathbf{z}_k(\boldsymbol{\beta}_{\text{dist},2}^{(t)}). \end{aligned}$$

Using straightforward algebraic calculations, we can show that, given $\boldsymbol{\alpha}$,

$$\begin{aligned} \mathcal{L}_N(\boldsymbol{\alpha}) &= \mathcal{L}_N(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) + \frac{1}{2} \text{vecl}(\boldsymbol{\alpha} - \boldsymbol{\beta}_{\text{dist},2}^{(t)})^T \widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) \text{vecl}(\boldsymbol{\alpha} - \boldsymbol{\beta}_{\text{dist},2}^{(t)}) \\ &\quad + \text{vecl}(\boldsymbol{\alpha} - \boldsymbol{\beta}_{\text{dist},2}^{(t)})^T \left\{ \widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) \text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) - \mathbf{z}_N(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) \right\}. \end{aligned}$$

There are three quantities on the right-hand side of the above. The first is irrelevant to $\boldsymbol{\alpha}$, and thus can be ignored in the optimization. The second involves $\widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},2}^{(t)})$, which is a $\{(p-d)d \times (p-d)d\}$ matrix. Transmitting $\widehat{\mathbf{H}}_k(\boldsymbol{\beta}_{\text{dist},2}^{(t)})$ from the local machines to the central one to form $\widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},2}^{(t)})$ is communication inefficient, particularly when the covariates are high dimensional. In parallel to the first distributed algorithm, we suggest replacing $\widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},2}^{(t)})$ with $\widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},2}^{(t)})$, to obtain

$$\begin{aligned} \widetilde{\mathcal{L}}_N(\boldsymbol{\alpha}) &\stackrel{\text{def}}{=} \mathcal{L}_N(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) + \frac{1}{2} \text{vecl}(\boldsymbol{\alpha} - \boldsymbol{\beta}_{\text{dist},2}^{(t)})^T \widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) \text{vecl}(\boldsymbol{\alpha} - \boldsymbol{\beta}_{\text{dist},2}^{(t)}) \\ &\quad + \text{vecl}(\boldsymbol{\alpha} - \boldsymbol{\beta}_{\text{dist},2}^{(t)})^T \left\{ \widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) \text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) - \mathbf{z}_N(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) \right\}. \end{aligned}$$

We perform optimization on the first machine. Using $\widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},2}^{(t)})$ in place of

$\widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},2}^{(t)})$ does not incur a communication cost. If $\boldsymbol{\alpha}$ is sufficiently close to $\boldsymbol{\beta}_{\text{dist},2}^{(t)}$, it is reasonable to expect the approximation error

$$\text{vecl}(\boldsymbol{\alpha} - \boldsymbol{\beta}_{\text{dist},2}^{(t)})^T \left\{ \widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) - \widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) \right\} \text{vecl}(\boldsymbol{\alpha} - \boldsymbol{\beta}_{\text{dist},2}^{(t)})$$

to be negligible. Next, we study the third quantity. Here, $\{\widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},2}^{(t)})\text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(t)})\}$ can be formed on each local machine using the relation

$$\widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},2}^{(t)})\text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) = m^{-1} \sum_{k=1}^m \{\widehat{\mathbf{H}}_k(\boldsymbol{\beta}_{\text{dist},2}^{(t)})\text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(t)})\}.$$

In particular, we form $\{\widehat{\mathbf{H}}_k(\boldsymbol{\beta}_{\text{dist},2}^{(t)})\text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(t)})\}$ and $\mathbf{z}_k(\boldsymbol{\beta}_{\text{dist},2}^{(t)})$ on each local machine, and transmit these random vectors to the central machine, incurring a communication cost of $O\{(p-d)d\}$, which is the minimal price that we have to pay for distributed algorithms. We ignore all quantities that are irrelevant to $\boldsymbol{\alpha}$ in $\widehat{\mathcal{L}}_N(\boldsymbol{\alpha})$. An equivalent form of $\widehat{\mathcal{L}}_N(\boldsymbol{\alpha})$ can be defined as

$$\begin{aligned} \mathcal{L}_N^*(\boldsymbol{\alpha}) &\stackrel{\text{def}}{=} \frac{1}{2} \text{vecl}(\boldsymbol{\alpha})^T \widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) \text{vecl}(\boldsymbol{\alpha}) \\ &\quad + \text{vecl}(\boldsymbol{\alpha})^T \left[\{\widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) - \widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},2}^{(t)})\} \text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) - \mathbf{z}_N(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) \right]. \end{aligned} \quad (4.2)$$

We seek $\boldsymbol{\beta}$ by minimizing $\mathcal{L}_N^*(\boldsymbol{\alpha})$ in distributed settings.

The second distributed algorithm proceeds as follows. We start from $\boldsymbol{\beta}_{\text{dist},2}^{(0)}$. Once we have $\boldsymbol{\beta}_{\text{dist},2}^{(t+1)}$ on the first machine, we update $\widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},2}^{(t)})$ with $\widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},2}^{(t+1)})$. We broadcast $\boldsymbol{\beta}_{\text{dist},2}^{(t+1)}$ from the first machine to the remaining machines to update $\widehat{\mathbf{H}}_k(\boldsymbol{\beta}_{\text{dist},2}^{(t+1)})\text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(t+1)}) - \mathbf{z}_k(\boldsymbol{\beta}_{\text{dist},2}^{(t+1)})$, which are then transmitted to the first machine to update the quantity in the square brackets in (4.2). Define

$$\begin{aligned} \boldsymbol{\beta}_{\text{dist},2}^{(t+2)} &\stackrel{\text{def}}{=} \underset{\boldsymbol{\alpha}}{\text{argmin}} \left(\frac{1}{2} \text{vecl}(\boldsymbol{\alpha})^T \widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},2}^{(t+1)}) \text{vecl}(\boldsymbol{\alpha}) \right. \\ &\quad \left. + \left[\{\widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},2}^{(t+1)}) - \widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},2}^{(t+1)})\} \text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(t+1)}) - \mathbf{z}_N(\boldsymbol{\beta}_{\text{dist},2}^{(t+1)}) \right] + \lambda \|\boldsymbol{\alpha}\|_1 \right). \end{aligned}$$

We iterate until convergence, and denote the final solution by $\widehat{\boldsymbol{\beta}}_{\text{dist},2}$.

4.2. Theoretical properties of the second distributed algorithm

Let \mathcal{S} be the support of $\text{vecl}(\boldsymbol{\beta})$, and the cardinality of \mathcal{S} be $s \stackrel{\text{def}}{=} |\mathcal{S}|$. We introduce conditions (C5')–(C8') to replace (C5)–(C8) to study the theoretical properties of $\widehat{\boldsymbol{\beta}}_{\text{dist},2}$.

(C5') The response Y is bounded, or the error ε has sub-Gaussian tails. The mean function satisfies $\sup |m(\mathbf{x}^T \boldsymbol{\beta})| < \infty$ and $E\{\|\mathbf{m}_1(\mathbf{x}^T \boldsymbol{\beta})\|_2^4\} < \infty$. The

loss $\mathcal{L}(\boldsymbol{\alpha})$ is restricted strongly convex over \mathcal{S} : for all $\delta \in C(\mathcal{S}) \stackrel{\text{def}}{=} \{\nu : |\nu_{\mathcal{S}}|_1 \leq 3|\nu_{\mathcal{S}^c}|_1, \mathcal{L}_1(\boldsymbol{\beta} + \delta) - \mathcal{L}_1(\boldsymbol{\beta}) - \delta \nabla \mathcal{L}_1(\boldsymbol{\beta}) \geq \mu |\delta|_2^2\}$. The Hessian is restricted Lipschitz: for all $\delta \in C(\mathcal{S})$, $|\{\widehat{\mathbf{H}}_1(\boldsymbol{\beta} + \delta) - \widehat{\mathbf{H}}_1(\boldsymbol{\beta})\}\delta|_{\infty} \leq M|\delta|_2^2$ and $|\{\widehat{\mathbf{H}}(\boldsymbol{\beta} + \delta) - \widehat{\mathbf{H}}(\boldsymbol{\beta})\}\delta|_{\infty} \leq M|\delta|_2^2$.

- (C6') The bandwidths satisfy $Nh_k^{2q}h_l^{2q} \rightarrow 0$, $Nh_1^{2(q-1)}h_l^{2q} \rightarrow 0$, $Nh_k^d h_l^d / \log^2 N \rightarrow \infty$ and $Nh_1^{d+2} / \log N \rightarrow \infty$. In addition, $s \log p = o\{N^{1/4} + h_k^{(-q+1)/2}h_l^{(-q+1)/2} + (Nh_1^{d+2} / \log N)^{1/2}\}$ for $1 \leq k \leq l \leq 4$.
- (C7') The covariate dimension p satisfies $p = O(N^{c_3}) \log p = O(N^{c_4})$, for $c_3 > 0$, $0 < c_4 < 1$. The sample size n satisfies $n = O(N^{c_5})$, for $0 < c_5 < 1$, and the sparsity level s satisfies $s = O(n^{c_6})$, for $0 \leq c_6 < 1/2$.
- (C8') The initial value $\boldsymbol{\beta}_{\text{dist},2}^{(0)}$ satisfies $|\text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(0)} - \boldsymbol{\beta})|_1 = O_p\{s(\log p/n)^{1/2}\}$ and $|\text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(0)} - \boldsymbol{\beta})|_2 = O_p(s \log p/n)^{1/2}$.

Theorem 2 provides an error bound for our distributed estimate $\boldsymbol{\beta}_{\text{dist},2}^{(t)}$.

Theorem 2. *Take*

$$\lambda_N^{(t)} = C \left\{ \left(\frac{\log p}{N} \right)^{1/2} + \left(\frac{\log p}{n} \right)^{1/2} |\text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(t-1)} - \boldsymbol{\beta})|_1 + |\text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(t-1)} - \boldsymbol{\beta})|_2^2 \right\},$$

for a sufficiently large constant C . Under conditions (C1)–(C4) and (C5')–(C8'), we have, for $t \geq 1$,

$$|\text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(t)} - \boldsymbol{\beta})|_2 = O_p \left\{ \left(s \frac{\log p}{N} \right)^{1/2} + s^{(2t+1)/2} \left(\frac{\log p}{n} \right)^{(t+1)/2} \right\}.$$

Theorem 2 indicates that, as the iteration proceeds, $\boldsymbol{\beta}_{\text{dist},2}^{(t)}$ improves accordingly. Indeed, $|\text{vecl}(\boldsymbol{\beta}_{\text{dist},2}^{(t)} - \boldsymbol{\beta})|_2$ is upper bounded by two orders: $O\{(s \log p/N)^{1/2}\}$ and $O\{s^{(2t+1)/2}(\log p/n)^{(t+1)/2}\}$. As long as the iteration step t is sufficiently large that $t \geq \log(p/n)/\log\{cn/(s^2 \log p)\}$, for some $c > 0$, the second order is dominated by the first, and the convergence rate of $\boldsymbol{\beta}_{\text{dist},2}^{(t)}$ becomes $O\{(s \log p/N)^{1/2}\}$, which is the oracle rate of $\widehat{\boldsymbol{\beta}}_{\text{pool},2}$. By condition (C7'), $\log(p/n)/\log\{cn/(s^2 \log p)\}$ is upper bounded, indicating that the difference between $\boldsymbol{\beta}_{\text{dist},2}^{(t)}$ and $\widehat{\boldsymbol{\beta}}_{\text{pool},2}$ is asymptotically negligible after a finite number of iterations. In other words, the distributed algorithm behaves asymptotically as well as the pooled algorithm.

Theorem 3. *In addition to conditions (C1)–(C4) and (C5')–(C8'), we further assume that $\|\boldsymbol{\Sigma}_{\mathcal{S}^c \times \mathcal{S}} \boldsymbol{\Sigma}_{\mathcal{S} \times \mathcal{S}}^{-1}\|_{\infty} \leq 1 - \alpha$, for some $0 < \alpha < 1$.*

1. *The distributed estimate satisfies $\mathcal{S}(\boldsymbol{\beta}_{\text{dist},2}^{(t)}) \subseteq \mathcal{S}$, with probability approaching one.*

2. Suppose for a sufficiently large constant C that

$$\min_{j \in \mathcal{S}} |\beta_j| \geq C \|\Sigma_{\mathcal{S} \times \mathcal{S}}^{-1}\|_{\infty} \left\{ \left(\frac{\log p}{N} \right)^{1/2} + s^t \left(\frac{\log p}{n} \right)^{(t+1)/2} \right\}.$$

Then, we have $\mathcal{S}(\beta_{\text{dist},2}^{(t)}) = \mathcal{S}$, with probability approaching one.

The irrepresentable condition and the “beta-min” condition are widely used in prior studies to establish the support recovery property; see, for example, Wainwright (2009). For $t \geq \log(p/n)/\log\{cn/(s^2 \log p)\}$, the “beta-min” condition reduces to $\min_{j \in \mathcal{S}} |\beta_j| \geq C \|\Sigma_{\mathcal{S} \times \mathcal{S}}^{-1}\|_{\infty} (\log p/N)^{1/2}$, which is a classic and widely used condition.

5. The Third Distributed Algorithm under Orthogonality Constraints

Let β be a basis of $\mathcal{S}_{E(Y|\mathbf{x})}$. In Sections 2–4, we required the upper $d \times d$ block of β to be $\mathbf{I}_{d \times d}$, which ensures that β is identifiable and that $\text{vecl}(\beta)$ are all free parameters. However, this implicitly requires that the first d covariates of \mathbf{x} be truly predictive, which is not always realistic. In this section, we merely assume β is orthonormal, such that $\beta^T \beta = \mathbf{I}_{d \times d}$. This orthogonality constraint does not ensure β is identifiable. In addition, optimization over the manifold $\{\beta : \beta^T \beta = \mathbf{I}_{d \times d}\}$ is, in general, difficult. However, this orthogonality constraint appears more realistic, and suffices to recover $\mathcal{S}_{E(Y|\mathbf{x})}$. We propose a third distributed algorithm under this orthogonality constraint, following Wen and Yin (2013).

This section requires a few changes in notation. In particular, we redefine

$$\hat{\mathbf{x}}_{i,k}(\beta_{\text{dist},3}^{(t)}) \stackrel{\text{def}}{=} \text{vec} \left\{ \left[\mathbf{x}_{i,k} - \frac{\hat{E}\{\mathbf{x}_{i,k} w(\mathbf{x}_{i,k}) \mid \mathbf{x}_{i,k}^T \beta_{\text{dist},3}^{(t)}\}}{\hat{E}\{w(\mathbf{x}_{i,k}) \mid \mathbf{x}_{i,k}^T \beta_{\text{dist},3}^{(t)}\}} \right] \hat{\mathbf{m}}_1^T(\mathbf{x}_{i,k}^T \beta_{\text{dist},3}^{(t)}) \right\},$$

where we use “vec” in place of “vecl” in (3.1). Accordingly,

$$\hat{Y}_{i,k}(\beta_{\text{dist},3}^{(t)}) \stackrel{\text{def}}{=} \{\hat{\mathbf{x}}_{i,k}^T(\beta_{\text{dist},3}^{(t)})\} \text{vec}(\beta_{\text{dist},3}^{(t)}) + \{Y_{i,k} - \hat{m}(\mathbf{x}_{i,k}^T \beta_{\text{dist},3}^{(t)})\},$$

where $\beta_{\text{dist},3}^{(t)}$ is an intermediate estimate of β at the t th step. Redefine

$$\begin{aligned} \mathbf{z}_k(\beta_{\text{dist},3}^{(t)}) &\stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n w(\mathbf{x}_{i,k}) \hat{\mathbf{x}}_{i,k}(\beta_{\text{dist},3}^{(t)}) \hat{Y}_{i,k}(\beta_{\text{dist},3}^{(t)}) \text{ and} \\ \mathbf{z}_N(\beta_{\text{dist},3}^{(t)}) &\stackrel{\text{def}}{=} \frac{1}{m} \sum_{k=1}^m \mathbf{z}_k(\beta_{\text{dist},3}^{(t)}). \end{aligned}$$

In addition,

$$\hat{\mathbf{H}}_k(\beta_{\text{dist},3}^{(t)}) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n w(\mathbf{x}_{i,k}) \hat{\mathbf{x}}_{i,k}(\beta_{\text{dist},3}^{(t)}) \hat{\mathbf{x}}_{i,k}^T(\beta_{\text{dist},3}^{(t)}) \text{ and}$$

$$\widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) \stackrel{\text{def}}{=} \frac{1}{m} \sum_{k=1}^m \widehat{\mathbf{H}}_k(\boldsymbol{\beta}_{\text{dist},3}^{(t)}).$$

In parallel to (4.2), we redefine the loss function using this notation as

$$\begin{aligned} \mathcal{L}_N^*(\boldsymbol{\alpha}) &\stackrel{\text{def}}{=} \frac{1}{2} \text{vec}(\boldsymbol{\alpha})^T \widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) \text{vec}(\boldsymbol{\alpha}) \\ &\quad + \text{vec}(\boldsymbol{\alpha})^T \left[\{\widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) - \widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},3}^{(t)})\} \text{vec}(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) - \mathbf{z}_N(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) \right]. \end{aligned} \quad (5.1)$$

In this section, we propose an orthogonality-constrained optimization approach within the Stiefel manifold (Edelman, Arias and Smith, 1998). In other words, we minimize $\mathcal{L}_N^*(\boldsymbol{\alpha})$ subject to the orthogonality constraint $\boldsymbol{\alpha}^T \boldsymbol{\alpha} = \mathbf{I}_{d \times d}$.

Wen and Yin (2013) propose a first order descent algorithm, that yields a feasible solution, in that it preserves the updates within the manifold. Specifically, we define the gradient of $\mathcal{L}_N^*(\boldsymbol{\alpha})$ with respect to $\boldsymbol{\alpha}$ as

$$\begin{aligned} \mathbf{G}(\boldsymbol{\alpha}) &\stackrel{\text{def}}{=} \text{vec}^{-1} \left[\widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) \text{vec}(\boldsymbol{\alpha}) \right. \\ &\quad \left. + \{\widehat{\mathbf{H}}(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) - \widehat{\mathbf{H}}_1(\boldsymbol{\beta}_{\text{dist},3}^{(t)})\} \text{vec}(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) - \mathbf{z}_N(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) \right]. \end{aligned}$$

Define $\mathbf{W}(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) \stackrel{\text{def}}{=} \mathbf{G}(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) (\boldsymbol{\beta}_{\text{dist},3}^{(t)})^T - (\boldsymbol{\beta}_{\text{dist},3}^{(t)}) \{\mathbf{G}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})\}^T$, which is a skew-symmetric matrix. Let $\tau^{(t)}$ be the step size. In practice, $\tau^{(t)}$ can be chosen using a non-monotone line search with the Barzilai–Borwein step size (Barzilai and Borwein, 1988). By the Cayley transformation, we have $\boldsymbol{\beta}_{\text{dist},3}^{(t+1)} \stackrel{\text{def}}{=}$

$$\left\{ \mathbf{I}_{p \times p} + \tau^{(t)} \frac{\mathbf{W}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})}{2} \right\}^{-1} \left\{ \mathbf{I}_{p \times p} - \tau^{(t)} \frac{\mathbf{W}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})}{2} \right\} \{\boldsymbol{\beta}_{\text{dist},3}^{(t)}}. \quad (5.2)$$

It can be verified that $\{\boldsymbol{\beta}_{\text{dist},3}^{(t+1)}\}^T \{\boldsymbol{\beta}_{\text{dist},3}^{(t+1)}\} = \mathbf{I}_{d \times d}$, if $\{\boldsymbol{\beta}_{\text{dist},3}^{(t)}\}^T \{\boldsymbol{\beta}_{\text{dist},3}^{(t)}\} = \mathbf{I}_{d \times d}$. Thus, this algorithm preserves the constraint exactly. Starting from $\boldsymbol{\beta}^{(0)}$, we iterate (5.2) until convergence. Denote the final solution by $\widehat{\boldsymbol{\beta}}_{\text{dist},3}$.

The inversion $\{\mathbf{I}_{p \times p} + \tau^{(t)} \mathbf{W}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})/2\}^{-1}$ dominates the computation for $\boldsymbol{\beta}_{\text{dist},3}^{(t+1)}$ in (5.2). However, we do not have to invert a $p \times p$ matrix. In particular, calculating this inversion is very cheap, because $\mathbf{W}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})$ is formed as the outer product of two low-rank matrices. Rewrite $\mathbf{W}(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) = \{\mathbf{U}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})\} \{\mathbf{V}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})\}^T$, for $\mathbf{U}(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) \stackrel{\text{def}}{=} [\mathbf{G}(\boldsymbol{\beta}_{\text{dist},3}^{(t)}), \boldsymbol{\beta}_{\text{dist},3}^{(t)}] \in \mathbb{R}^{p \times 2d}$ and $\mathbf{V}(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) \stackrel{\text{def}}{=} [\boldsymbol{\beta}_{\text{dist},3}^{(t)}, -\mathbf{G}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})] \in \mathbb{R}^{p \times 2d}$. As long as $\{\mathbf{I}_{2d \times 2d} + \tau^{(t)} \mathbf{V}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})^T \mathbf{U}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})/2\}$ is invertible, which is often the case, by Lemma 4 of Wen and Yin (2013), an equivalent form of (5.2) is $\boldsymbol{\beta}_{\text{dist},3}^{(t+1)} = \boldsymbol{\beta}_{\text{dist},3}^{(t)} - \tau^{(t)} \mathbf{U}(\boldsymbol{\beta}_{\text{dist},3}^{(t)}) \{\mathbf{I}_{(2d) \times (2d)} + \tau^{(t)} \mathbf{V}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})^T \mathbf{U}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})/2\}^{-1} \{\mathbf{V}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})\}^T \{\boldsymbol{\beta}_{\text{dist},3}^{(t)}\}$. By the very purpose of mean dimension reduction, d is far less than p . It is thus natural to expect that inverting $\{\mathbf{I}_{(2d) \times (2d)} + \tau^{(t)} \mathbf{V}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})^T \mathbf{U}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})/2\} \in \mathbb{R}^{2d \times 2d}$ is much easier than inverting

$\{\mathbf{I}_{p \times p} + \tau^{(t)} \mathbf{W}(\boldsymbol{\beta}_{\text{dist},3}^{(t)})/2\} \in \mathbb{R}^{p \times p}$. In this sense, our distributed algorithm with an orthogonality constraint is computationally efficient.

6. Simulation Studies

We conduct simulation studies to demonstrate the finite sample performance of our proposed distributed algorithms. We generate the observations from the following examples.

Example 1. We generate \mathbf{x} from a multivariate normal distribution with mean zero and covariance $\boldsymbol{\Sigma} = (0.5^{|i-j|})_{p \times p}$. We generate Y from a normal distribution with mean $m(\mathbf{x}^T \boldsymbol{\beta}) = \sin(2\mathbf{x}^T \boldsymbol{\beta}) + 2\exp(2 + \mathbf{x}^T \boldsymbol{\beta})$ and variance $\sigma^2(\mathbf{x}) = \log\{2 + (\mathbf{x}^T \boldsymbol{\beta})\}$. The first four components of $\boldsymbol{\beta}$ are $(1, 1, -1, 1)^T$, and all other entries are identically zero. In this example, $p = 16$ and $d = 1$.

Example 2. We generate \mathbf{x} independently from a uniform distribution defined on $[-2, 2]$. We generate Y from a normal distribution with mean $m(\mathbf{x}^T \boldsymbol{\beta}) = (\mathbf{x}^T \boldsymbol{\beta}_1)/\{0.5 + (1.5 + \mathbf{x}^T \boldsymbol{\beta}_2)^2\}$ and variance $\sigma^2(\mathbf{x}) = \exp(X_1)$, where X_1 is the first coordinate of \mathbf{x} . The first four components of $\boldsymbol{\beta}_1$ and $\boldsymbol{\beta}_2$ are $\boldsymbol{\beta}_1 = (1, 0, 1, 1)^T$ and $\boldsymbol{\beta}_2 = (0, 1, -1, 1)^T$, respectively. All other entries of $\boldsymbol{\beta}_1$ and $\boldsymbol{\beta}_2$ are identically zero. In this example, $p = 16$ and $d = 2$.

We run 500 replicates to compare the performance of the following estimates:

1. $\hat{\boldsymbol{\beta}}_{\text{pool},1}(w)$: The pooled estimate that pools all observations together and uses the true weight $w(\mathbf{x}) = \{\sigma^2(\mathbf{x})\}^{-1}$. This serves as a benchmark for algorithm 1.
2. $\hat{\boldsymbol{\beta}}_{\text{dist},1}(w)$: The distributed estimate that uses $w(\mathbf{x}) = \{\sigma^2(\mathbf{x})\}^{-1}$.
3. $\hat{\boldsymbol{\beta}}_{\text{dist},1}(w^*)$: The distributed estimate that uses $w^*(\mathbf{x}) = 1$.
4. $\hat{\boldsymbol{\beta}}_{\text{pool},2}(w)$: The regularized pooled estimate that aggregates all observations together and uses the true weight $w(\mathbf{x}) = \{\sigma^2(\mathbf{x})\}^{-1}$. This serves as a benchmark for algorithm 2.
5. $\hat{\boldsymbol{\beta}}_{\text{dist},2}(w)$: The regularized distributed estimate that uses $w(\mathbf{x}) = \{\sigma^2(\mathbf{x})\}^{-1}$.
6. $\hat{\boldsymbol{\beta}}_{\text{dist},2}(w^*)$: The regularized distributed estimate that misspecifies $w(\mathbf{x})$ as $w^*(\mathbf{x}) = 1$.
7. $\hat{\boldsymbol{\beta}}_{\text{pool},3}(w)$: The pooled estimate that aggregates all observations together and uses the true weight $w(\mathbf{x}) = \{\sigma^2(\mathbf{x})\}^{-1}$. This serves as a benchmark for algorithm 3.
8. $\hat{\boldsymbol{\beta}}_{\text{dist},3}(w)$: The distributed estimate that uses $w(\mathbf{x}) = \{\sigma^2(\mathbf{x})\}^{-1}$.
9. $\hat{\boldsymbol{\beta}}_{\text{dist},3}(w^*)$: The distributed estimate that misspecifies $w(\mathbf{x})$ as $w^*(\mathbf{x}) = 1$.

Let β be a basis matrix of the central mean subspace, and $\widehat{\beta}$ be its estimate. To assess the estimation accuracy of $\widehat{\beta}$, we use the Euclidean distance between β and $\widehat{\beta}$, defined as the Frobenius norm of the matrix $\widehat{\beta}(\widehat{\beta}^T\widehat{\beta})^{-1}\widehat{\beta}^T - \beta(\beta^T\beta)^{-1}\beta^T$. A smaller distance indicates a better estimate.

Throughout, we fix the total sample size $N = 2,500$. We consider three combinations, $(n, m) = (500, 5)$, $(250, 10)$, and $(100, 25)$, where m is the number of machines. We choose the initial value $\beta^{(0)}$ for algorithms 1 and 3 by using a minimum average variance estimation (Xia et al., 2002). For the initial value of algorithm 2, we implement the sparse sliced inverse regression (Lin, Zhao and Liu, 2019). We choose the bandwidths using a “rule-of-thumb” approach because the semiparametric estimating equations approach is not sensitive to the bandwidth selections (Ma and Zhu, 2014). In particular, we set $h_1 = h_2 = h_3 = cn^{-1/(4+d)}$.

Tables 1 summarizes the average distances and CPU running times (in seconds) of various estimates. For $j \in \{1, 2, 3\}$, the pooled estimate $\widehat{\beta}_{\text{pool},j}(w)$ performs best, in that it has the smallest biases across all scenarios. Furthermore, the biases of the distributed estimates, $\widehat{\beta}_{\text{dist},j}(w)$ and $\widehat{\beta}_{\text{dist},j}(w^*)$, increase with the number of machines. Not surprisingly, $\widehat{\beta}_{\text{dist},j}(w^*)$ is relatively less accurate among the distributed estimates because the weight function is misspecified. The distributed algorithms reduce the computational complexity substantially. Algorithm 1 is slightly faster, but less accurate than algorithm 3. Algorithm 2 has the smallest distance of the distributed algorithms, because it benefits from a sparse structure. However, it requires the most computational resources.

7. Conclusion

In this paper, we have introduced distributed algorithms for estimating the central mean subspace under two sets of identifiability conditions. The first set requires that the upper block of the basis of the central mean subspace is an identity matrix. Under this condition, we design two distributed algorithms. The first produces a dense solution, which suffices if the covariate dimension is moderate. The second generates a sparse solution, which allows the covariates to be high or even ultrahigh dimensional. For the second distributed algorithm, an important contribution is that we recast the problems of estimating equations under a least squares framework. This enables us to incorporate an appropriate penalty to produce a solution, and more importantly, allows us to solve the penalized algorithms under a linear regression framework. This idea is interesting and can be adapted to solve other problems that use estimating equations.

The second set of identifiability conditions assumes that the basis of the central mean subspace is orthonormal. Here, determining a feasible solution is challenging. We address this problem using the first order descent algorithm. However, finding a sparse feasible solution remains challenging owing to the discontinuity of the sparse solution path, and thus warrants further research.

Table 1. The average distance and CPU running time (in seconds) of various estimates.

(n, m)	(500,5)		(250, 10)		(100, 25)	
	distance	time	distance	time	distance	time
Example 1						
$\hat{\beta}_{\text{pool},1}(w)$			0.224	(26.564)		
$\hat{\beta}_{\text{dist},1}(w)$	0.236	(10.391)	0.245	(7.613)	0.261	(3.158)
$\hat{\beta}_{\text{dist},1}(w^*)$	0.249	(8.484)	0.252	(5.322)	0.273	(2.613)
$\hat{\beta}_{\text{pool},2}(w)$			0.144	(33.216)		
$\hat{\beta}_{\text{dist},2}(w)$	0.157	(15.729)	0.169	(10.048)	0.178	(6.165)
$\hat{\beta}_{\text{dist},2}(w^*)$	0.163	(13.963)	0.184	(8.371)	0.197	(4.687)
$\hat{\beta}_{\text{pool},3}(w)$			0.219	(31.854)		
$\hat{\beta}_{\text{dist},3}(w)$	0.226	(13.432)	0.238	(9.583)	0.255	(5.227)
$\hat{\beta}_{\text{dist},3}(w^*)$	0.234	(10.038)	0.243	(7.255)	0.264	(3.641)
Example 2						
$\hat{\beta}_{\text{pool},1}(w)$			0.304	(31.859)		
$\hat{\beta}_{\text{dist},1}(w)$	0.315	(14.337)	0.322	(10.741)	0.341	(4.845)
$\hat{\beta}_{\text{dist},1}(w^*)$	0.323	(12.148)	0.334	(8.066)	0.358	(3.275)
$\hat{\beta}_{\text{pool},2}(w)$			0.169	(34.738)		
$\hat{\beta}_{\text{dist},2}(w)$	0.177	(18.344)	0.182	(13.765)	0.189	(7.148)
$\hat{\beta}_{\text{dist},2}(w^*)$	0.188	(16.731)	0.194	(11.142)	0.205	(5.671)
$\hat{\beta}_{\text{pool},3}(w)$			0.291	(32.530)		
$\hat{\beta}_{\text{dist},3}(w)$	0.299	(15.138)	0.316	(11.435)	0.322	(5.836)
$\hat{\beta}_{\text{dist},3}(w^*)$	0.314	(13.256)	0.327	(8.997)	0.349	(4.158)

Numerous works have proposed solutions to the problem of high dimensionality. Here, we focus on massive data of high dimensions and large volumes. We propose several distributed algorithms, which have nearly minimal communication cost and almost the lowest computational complexity. In addition, our solutions possess many desirable theoretical properties. However, adapting these distributed algorithms to identify and recover the central subspaces, particularly when the response variables are multivariate or even high dimensional, remains an open problem, and thus is left to future research.

Supplementary Material

The online Supplementary Material contains descriptions of the pooled algorithms, additional simulations, and technical proofs of all theorems.

Acknowledgments

This work was supported by the National Natural Science Foundation of China (12171477, 12225113) and Beijing Natural Science Foundation (Z190002). All correspondence can be directed to Liping Zhu.

References

Balcan, M. F., Liang, Y., Song, L., Woodruff, D. and Xie, B. (2016). Communication efficient distributed kernel principal component analysis. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 725–734.

Barzilai, J. and Borwein, J. M. (1988). Two-point step size gradient methods. *Journal of Numerical Analysis* **8**, 141–148.

Battey, H., Fan, J., Liu, H., Lu, J. and Zhu, Z. (2015). Distributed estimation and inference with statistical guarantees. *arXiv:1509.05457*.

Battey, H., Fan, J., Liu, H., Lu, J. and Zhu, Z. (2018). Distributed testing and estimation under sparse high dimensional models. *The Annals of Statistics* **46**, 1352–1382.

Cai, Z., Li, R. and Zhu, L. (2020). Online sufficient dimension reduction through sliced inverse regression. *Journal of Machine Learning Research* **21**, 1–25.

Chen, C., Xu, W. and Zhu, L. (2022). Distributed estimation in heterogeneous reduced rank regression: With application to order determination in sufficient dimension reduction. *Journal of Multivariate Analysis* **190**, 104991.

Cook, R. D. (2009). *Regression Graphics: Ideas for Studying Regressions through Graphics*. John Wiley & Sons.

Cook, R. D. and Li, B. (2002). Dimension reduction for conditional mean in regression. *The Annals of Statistics* **30**, 455–474.

Cook, R. D. and Li, B. (2004). Determining the dimension of iterative Hessian transformation. *The Annals of Statistics* **32**, 2501–2531.

Cook, R. D. and Weisberg, S. (1991). Discussion of sliced inverse regression for dimension reduction. *Journal of the American Statistical Association* **86**, 328–332.

Edelman, A., Arias, T. A. and Smith, S. T. (1998). The geometry of algorithms with orthogonality constraints. *SIAM Journal on Matrix Analysis and Applications* **20**, 303–353.

Fan, J., Guo, Y. and Wang, K. (2023). Communication-efficient accurate statistical estimation. *Journal of the American Statistical Association* **542**, 1000–1010.

Fan, J., Wang, D., Wang, K. and Zhu, Z. (2019). Distributed estimation of principal eigenspaces. *The Annals of Statistics* **47**, 3009–3031.

Härdle, W. and Stoker, T. M. (1989). Investigating smooth multiple regression by the method of average derivatives. *Journal of the American Statistical Association* **84**, 986–995.

Jordan, M. I., Lee, J. D. and Yang, Y. (2019). Communication-efficient distributed statistical inference. *Journal of the American Statistical Association* **114**, 668–681.

Li, K.-C. (1991). Sliced inverse regression for dimension reduction. *Journal of the American Statistical Association* **86**, 316–327.

Li, K.-C. (1992). On principal Hessian directions for data visualization and dimension reduction: Another application of Stein’s lemma. *Journal of the American Statistical Association* **87**, 1025–1039.

Li, K.-C. and Duan, N. (1989). Regression analysis under link violation. *The Annals of Statistics* **17**, 1009–1052.

Lin, Q., Zhao, Z. and Liu, J. S. (2019). Sparse sliced inverse regression via Lasso. *Journal of the American Statistical Association* **114**, 1726–1739.

Luo, W. and Cai, X. (2016). A new estimator for efficient dimension reduction in regression. *Journal of Multivariate Analysis* **145**, 236–249.

Luo, W., Li, B. and Yin, X. (2014). On efficient dimension reduction with respect to a statistical functional of interest. *The Annals of Statistics* **42**, 382–412.

Ma, Y. and Zhu, L. (2013). Efficiency loss and the linearity condition in dimension reduction. *Biometrika* **100**, 371–383.

Ma, Y. and Zhu, L. (2014). On estimation efficiency of the central mean subspace. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* **76**, 885–901.

Tibshirani, R. (1996). Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* **58**, 267–288.

Wainwright, M. J. (2009). Sharp thresholds for high-dimensional and noisy sparsity recovery using ℓ_1 -constrained quadratic programming (Lasso). *IEEE Transactions on Information Theory* **55**, 2183–2202.

Wen, Z. and Yin, W. (2013). A feasible method for optimization with orthogonality constraints. *Mathematical Programming* **142**, 397–434.

Xia, Y., Tong, H., Li, W. K. and Zhu, L.-X. (2002). An adaptive estimation of dimension reduction space. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* **64**, 363–410.

Zhang, Y., Duchi, J. C. and Wainwright, M. J. (2013). Communication-efficient algorithms for statistical optimization. *Journal of Machine Learning Research* **14**, 3321–3363.

Zhu, L. and Zhong, W. (2015). Estimation and inference on central mean subspace for multivariate response data. *Computational Statistics & Data Analysis* **92**, 68–83.

Zhu, Z. and Zhu, L. (2022). Distributed dimension reduction with nearly oracle rate. *Statistical Analysis and Data Mining* **15**, 692–706.

(Received May 2022; accepted April 2023)