



Published in final edited form as:

Ann Stat. 2018 June ; 46(3): 1352–1382. doi:10.1214/17-AOS1587.

DISTRIBUTED TESTING AND ESTIMATION UNDER SPARSE HIGH DIMENSIONAL MODELS

Heather Battey^{§,¶,*}, Jianqing Fan^{§,||,‡}, Han Liu[§], Junwei Lu[§], and Ziwei Zhu[§]

[§]Princeton University

[¶]Imperial College London

^{||}Fudan University

Abstract

This paper studies hypothesis testing and parameter estimation in the context of the divide-and-conquer algorithm. In a unified likelihood based framework, we propose new test statistics and point estimators obtained by aggregating various statistics from k subsamples of size n/k , where n is the sample size. In both low dimensional and sparse high dimensional settings, we address the important question of how large k can be, as n grows large, such that the loss of efficiency due to the divide-and-conquer algorithm is negligible. In other words, the resulting estimators have the same inferential efficiencies and estimation rates as an oracle with access to the full sample. Thorough numerical results are provided to back up the theory.

Keywords and phrases

Divide and conquer; debiasing; massive data; thresholding

MSC 2010 subject classifications

Primary 62F05; 62F10; secondary 62F12

1. Introduction

In recent years, the field of statistics has developed apace in response to the opportunities and challenges spawned from the ‘data revolution’, which marked the dawn of an era characterized by the availability of enormous datasets. An extensive toolkit of methodology is now in place for addressing a wide range of high dimensional problems, whereby the number of unknown parameters, d , is much larger than the number of observations, n . However, many modern datasets are instead characterized by n and d both large. The latter presents intimidating practical challenges resulting from storage and computational limitations, as well as numerous statistical challenges (Fan et al., 2014). It is important that

*Department of Mathematics, Imperial College London, London, SW7 2AZH, UK. Heather Battey was supported in part by the NIH grant 2R01-GM072611-11 as a postdoctoral fellow at Princeton University.

[§]Department of Operations Research and Financial Engineering, Princeton University, Princeton, NJ 08540

[‡]Jianqing Fan was supported in part by NSF Grants DMS-1206464 and DMS-1406266, and NIH grants 2R01-GM072611-11.

statistical methodology targeting modern application areas does not lose sight of the practical burdens associated with manipulating such large scale datasets. In this vein, incisive new algorithms have been developed for exploiting modern computing architectures and recent advances in distributed computing. These algorithms enjoy computational or communication efficiency and facilitate data handling and storage, but come with a statistical overhead if inappropriately tuned.

With increased mindfulness of the algorithmic difficulties associated with large datasets, the statistical community has witnessed a surge in recent activity in the statistical analysis of various divide and conquer (DC) algorithms, which randomly partition the n observations into k subsamples of size $n_k = n/k$, construct statistics based on each subsample, and aggregate them in a suitable way. In splitting the dataset, a single, very large scale estimation or testing problem with computational complexity $\mathcal{O}(\chi(n))$, for a given function $\chi(\cdot)$ that depends on the underlying problem, is transformed into k smaller problems with computational complexity $\mathcal{O}(\chi(n/k))$ on each machine. What get lost in this process are the interactions of split subsamples in each machine. They are not recoverable without additional rounds of communication or without additional communication between the machines. Since every additional split of the dataset incurs some efficiency loss, it is of significant practical interest to derive a theoretical upper bound on the number of subsamples k that delivers the same asymptotic statistical performance as the practically unavailable “oracle” procedure based on the full sample.

We develop communication efficient generalizations of the Wald and Rao’s score tests for the sparse high dimensional scheme, as well as communication efficient estimators for the parameters of the sparse high dimensional and low dimensional linear and generalized linear models. In all cases we give the upper bound on k for preserving the statistical error of the analogous full sample procedure. While hypothesis testing in a low dimensional context is straightforward, in the sparse high dimensional setting, nuisance parameters introduce a non-negligible bias, causing classical low dimensional theory to break down. In our high dimensional Wald construction, the phenomenon is remedied through a debiasing of the estimator, which gives rise to a test statistic with tractable limiting distribution, as documented in the $k = 1$ (no sample split) setting in Zhang and Zhang (2014) and van de Geer et al. (2014). For the high dimensional analogue of Rao’s score statistic, the incorporation of a correction factor increases the convergence rate of higher order terms, thereby vanquishing the effect of the nuisance parameters. The approach is introduced in the $k = 1$ setting in Ning and Liu (2014), where the test statistic is shown to possess a tractable limit distribution. However, the computation complexity for the debiased estimators increases by an order of magnitude, due to solving d high-dimensional regularization problems. This motivates us to appeal to the divide and conquer strategy.

We develop the theory and methodology for DC versions of these tests. In the case of $k = 1$, each of the above test statistics can be decomposed into a dominant term with tractable limit distribution and a negligible remainder term. The DC extension requires delicate control of these remainder terms to ensure the error accumulation remains sufficiently small so as not to materially contaminate the leading term. We obtain an upper bound on the number of permitted subsamples, k , subject to a statistical guarantee. More specifically, we find that the

theoretical upper bound on the number of subsamples guaranteeing the same inferential or estimation efficiency as the whole-sample procedure is $k = o((s \log d)^{-1} \sqrt{n})$ in the linear model, where s is the sparsity of the parameter vector. In the generalized linear model the scaling is $k = o(((s \vee s_1) \log d)^{-1} \sqrt{n})$, where s_1 is the sparsity of the inverse information matrix.

For sparse high dimensional estimation problems, we use the same debiasing technique introduced in the high dimensional testing problems to obtain a thresholded divide and conquer estimator that achieves the full sample minimax rate. The appropriate scaling is found to be $k = O(\sqrt{n/(s^2 \log d)})$ for the estimation of the sparse parameter vector in the high dimensional linear model and $k = O(\sqrt{n/((s \vee s_1)^2 \log d)})$ for the high dimensional generalized linear model. Moreover, we find that the loss incurred by the divide and conquer strategy, as quantified by the distance between the DC estimator and the full sample estimator, is negligible in comparison to the statistical error of the full sample estimator provided that k is not too large. In the context of estimation, the optimal scaling of k with n and d is also developed for the low dimensional linear and generalized linear model. This theory is of independent interest. It also allows us to study a refitted estimation procedure under a minimal signal strength assumption.

1.1. Related Literature

A partial list of references covering DC algorithms from a statistical perspective is Chen and Xie (2012), Zhang et al. (2013), Kleiner et al. (2014), Liu and Ihler (2014) and Zhao et al. (2014a). The closest works to ours are Zhang et al. (2013), Lee et al. (2015) and Rosenblatt and Nadler (2016). Zhang et al. (2013) consider the distributed estimator for kernel ridge regression. In the context of $d < n$, Zhang et al. (2013) propose the distributed estimator by averaging the kernel ridge regression estimators for each data split. They obtain an explicit upper bound on the number of splits yielding the minimax optimal rates for the mean squared error. However, it is not straightforward to generalize their estimator to the high dimensional setting. In an independent work, Lee et al. (2015) propose the same debiasing approach of van de Geer et al. (2014) to allow aggregation of local estimates on distributed data splits in the context of sparse high dimensional linear and generalized linear models. Though using different techniques of proofs, the conclusions of Lee et al. (2015) in terms of the optimal choice of tuning parameter scaling and the upper bound on the permissible number of sample splits is of the same order. Our work differs from theirs in two aspects: (1) our work also contributes to the distributed testing in sparse high dimensional models and (2) we propose a refitted distributed estimator which has the oracle rate. Our results on hypothesis testing reveal a different phenomenon to that found in estimation, as we observe through the different requirements on the scaling of k . On the estimation side, our results also differ from those of Lee et al. (2015) in that our additional refitting step allows us to achieve the oracle rate. Rosenblatt and Nadler (2016) consider the distributed empirical risk minimization for M -estimators. They require the dimension of the interest parameter to satisfy the scaling condition $d/n \rightarrow \kappa \in (0, 1)$, which rules out the $d \gg n$ case. They quantify the accuracy loss over the full sample estimator in terms of the number of splits.

1.2. Organization of the paper

The rest of the paper is organized as follows. Section 2 collects notation and details of a generic likelihood based framework. Section 3 covers testing, providing high dimensional DC analogues of the Wald test (Section 3.1) and Rao’s score test (Section 3.2), in each case deriving a tractable limit distribution for the corresponding test statistic under standard assumptions. Section 4 covers distributed estimation, proposing an aggregated estimator of the unknown parameters of linear and generalized linear models in low dimensional and sparse high dimensional scenarios, as well as a refitting procedure that improves the estimation rate, with the same scaling, under a minimal signal strength assumption. Section 5 provides numerical experiments to back up the developed theory. In Section 6 we discuss our results together with remaining future challenges. Proofs of our main results are collected in Section 7, while the statement and proofs of a number of technical lemmas are deferred to the Supplementary Material.

2. Background and Notation

We first collect the general notation, before providing a formal statement of our statistical problems. More specialized notation is introduced in context.

2.1. Generic Notation

We adopt the common convention of using bold-face letters for vectors only, while regular font is used for both matrices and scalars. $| \cdot |$ denotes both absolute value and cardinality of a set, with the context ensuring no ambiguity. For $\mathbf{x} = (x_1, \dots, x_d)^T \in \mathbb{R}^d$, and $1 \leq q < \infty$, we define $\|\mathbf{x}\|_q = (\sum_{j=1}^d |x_j|^q)^{1/q}$, $\|\mathbf{x}\|_0 = |\text{supp}(\mathbf{x})|$, where $\text{supp}(\mathbf{x}) = \{j : x_j \neq 0\}$. Write $\|\mathbf{x}\|_\infty = \max_{j=1, \dots, d} |x_j|$, while for a matrix $M = [M_{jk}]$, let $\|M\|_{\max} = \max_{j,k} |M_{jk}|$, $\|M\|_1 = \sum_{j,k} |M_{jk}|$. For any matrix M we use M_{ℓ} to index the transposed ℓ^{th} row of M and $[M]_{\ell}$ to index the ℓ^{th} column. The sub-Gaussian norm of a scalar random variable X is defined as $\|X\|_{\psi_2} = \sup_{q \geq 1} q^{-1/2} (\mathbb{E}|X|^q)^{1/q}$. For a random vector $X \in \mathbb{R}^d$, its sub-Gaussian norm is defined as $\|X\|_{\psi_2} = \sup_{\mathbf{x} \in \mathbb{S}^{d-1}} \|\langle X, \mathbf{x} \rangle\|_{\psi_2}$, where \mathbb{S}^{d-1} denotes the unit sphere in \mathbb{R}^d . Let I_d denote the $d \times d$ identity matrix; when the dimension is clear from the context, we omit the subscript. We also denote the Hadamard product of two matrices A and B as $A \circ B$ and $(A \circ B)_{jk} = A_{jk} B_{jk}$ for any j, k . $\{e_1, \dots, e_d\}$ denotes the canonical basis for \mathbb{R}^d . For a vector $\mathbf{v} \in \mathbb{R}^d$ and a set of indices $\mathcal{S} \subseteq \{1, \dots, d\}$, $\mathbf{v}_{\mathcal{S}}$ is the vector of length $|\mathcal{S}|$ whose components are $\{v_j : j \in \mathcal{S}\}$. Additionally, for a vector \mathbf{v} with j^{th} element v_j , we use the notation \mathbf{v}_{-j} to denote the remaining vector when the j^{th} element is removed. With slight abuse of notation, we write $\mathbf{v} = (v_j, \mathbf{v}_{-j})$ when we wish to emphasize the dependence of \mathbf{v} on v_j and \mathbf{v}_{-j} individually. The gradient of a function $f(\mathbf{x})$ is denoted by $\nabla f(\mathbf{x})$, while $\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})$ denotes the gradient of $f(\mathbf{x}, \mathbf{y})$ with respect to \mathbf{x} , and $\nabla_{\mathbf{xy}}^2 f(\mathbf{x}, \mathbf{y})$ denotes the matrix of cross partial derivatives with respect to the elements of \mathbf{x} and \mathbf{y} . For a scalar η , we simply write $f'(\eta) := \nabla_{\eta} f(\eta)$ and $f''(\eta) := \nabla_{\eta\eta}^2 f(\eta)$. For a random variable X and a sequence of random variables, $\{X_n\}$, we write $X_n \rightsquigarrow X$ when $\{X_n\}$ converges weakly to X . If X is a random variable with standard distribution, say F_X , we simply write $X_n \rightsquigarrow F_X$. Given $a, b \in \mathbb{R}$, let $a \vee b$ and $a \wedge b$ denote

the maximum and minimum of a and b . We also make use of the notation $a_n \lesssim b_n$ ($a_n \gtrsim b_n$) if a_n is less than (greater than) b_n up to a constant, and $a_n \asymp b_n$ if a_n is the same order as b_n .

2.2. General Likelihood based Framework

Let $(X_1^T, Y_1^T)^T, \dots, (X_n^T, Y_n^T)^T$ be n i.i.d. copies of the random vector $(X^T, Y)^T$, whose realizations take values in $\mathbb{R}^d \times \mathcal{Y}$. Write the collection of these n i.i.d. random couples as $\mathcal{D} = \{(X_1^T, Y_1^T)^T, \dots, (X_n^T, Y_n^T)^T\}$ with $Y = (Y_1, \dots, Y_n)^T$ and $X = (X_1, \dots, X_n)^T \in \mathbb{R}^{n \times d}$.

Conditional on X_i , we assume Y_i is distributed as F_{β^*} for all $i \in \{1, \dots, n\}$, where F_{β^*} is a known distribution parameterized by a sparse d -dimensional vector β^* and has a density or mass function f_{β^*} . We thus define the negative log-likelihood function, $\ell_n(\beta)$, as

$$\ell_n(\beta) = \frac{1}{n} \sum_{i=1}^n \ell_i(\beta) = -\frac{1}{n} \sum_{i=1}^n \log f_{\beta}(Y_i | X_i). \quad (2.1)$$

We use $J^* = J(\beta^*)$ to denote the information matrix and Θ^* to denote $(J^*)^{-1}$, where $J(\beta) = \mathbb{E}[\nabla_{\beta}^2 \ell_n(\beta)]$.

For testing problems, our goal is to test $H_0: \beta_v^* = \beta_v^H$ for a specific fixed index $v \in \{1, \dots, d\}$.

We partition β^* as $\beta^* = (\beta_v^*, \beta_{-v}^{*T})^T \in \mathbb{R}^d$, where $\beta_{-v}^* \in \mathbb{R}^{d-1}$ is a vector of nuisance parameters and β_v^* is the parameter of interest. To handle the curse of dimensionality, we exploit a penalized M-estimator defined as,

$$\hat{\beta}^\lambda = \underset{\beta}{\operatorname{argmin}} \{ \ell_n(\beta) + \mathcal{P}_\lambda(\beta) \}, \quad (2.2)$$

with $\mathcal{P}_\lambda(\beta)$ a sparsity inducing penalty function with a regularization parameter λ . Examples of $\mathcal{P}_\lambda(\beta)$ include the convex ℓ_1 penalty, $\mathcal{P}_\lambda(\beta) = \lambda \|\beta\|_1 = \lambda \sum_{j=1}^d |\beta_j|$ which, in the context of the linear model, gives rise to the Lasso estimator (Tibshirani, 1996),

$$\hat{\beta}_{\text{Lasso}}^\lambda = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2n} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}. \quad (2.3)$$

Other penalties include folded concave penalties such as the smoothly clipped absolute deviation (SCAD) penalty (Fan and Li, 2001) and minimax concave MCP penalty (Zhang, 2010), which eliminate the estimation bias and attain the oracle rates of convergence (Loh and Wainwright, 2013; Wang et al., 2014a). The SCAD penalty is defined as

$$\mathcal{P}_\lambda(\boldsymbol{\beta}) = \sum_{j=1}^d p_\lambda(\beta_j), \text{ where } p_\lambda(t) = \int_0^{|t|} \left\{ \lambda \mathbb{1}(z \leq \lambda) + \frac{a\lambda - z}{a-1} \mathbb{1}(z > \lambda) \right\} dz, \quad (2.4)$$

for a given parameter $a > 0$ and MCP penalty is given by

$$\mathcal{P}_\lambda(\boldsymbol{\beta}) = \sum_{j=1}^d p_\lambda(\beta_j), \text{ where } p_\lambda(t) = \lambda \int_0^{|t|} \left(1 - \frac{z}{\lambda b}\right)_+ dz \quad (2.5)$$

where $b > 0$ is a fixed parameter. The only requirement we have on $P_\lambda(\boldsymbol{\beta})$ is that it induces an estimator satisfying the following condition.

Condition 2.1: For any $\delta \in (0, 1)$, if $\lambda \asymp \sqrt{\log(d/\delta)/n}$,

$$\mathbb{P}\left(\|\hat{\boldsymbol{\beta}}^\lambda - \boldsymbol{\beta}^*\|_1 > Csn^{-1/2}\sqrt{\log(d/\delta)}\right) \leq \delta, \quad (2.6)$$

where s is the sparsity of $\boldsymbol{\beta}^*$, i.e., $s = \|\boldsymbol{\beta}^*\|_0$.

Condition 2.1 is crucial for the theory developed in Sections 3 and 4. Under suitable conditions on the design matrix X , it holds for the Lasso, SCAD and MCP. See Bühlmann and van de Geer (2011); Fan and Li (2001); Zhang (2010) respectively and Zhang and Zhang (2012).

The DC algorithm randomly and evenly partitions \mathcal{D} into k disjoint subsets $\mathcal{D}_1, \dots, \mathcal{D}_k$, so that $\mathcal{D} = \cup_{j=1}^k \mathcal{D}_j$, $\mathcal{D}_j \cap \mathcal{D}_\ell = \emptyset$ for all $j, \ell \in \{1, \dots, k\}$, and $|\mathcal{D}_1| = |\mathcal{D}_2| = \dots = |\mathcal{D}_k| = n_k = n/k$, where it is implicitly assumed that n can be divided evenly. Let $\mathcal{Q}_j \subset \{1, \dots, n\}$ be the index set corresponding to the elements of \mathcal{D}_j . Then for an arbitrary $n \times d$ matrix A , $A^{(j)} = [A_{i\ell}]_{i \in \mathcal{Q}_j, \ell \in \mathcal{D}}$. For an arbitrary estimator $\hat{\boldsymbol{\tau}}$, we write $\hat{\boldsymbol{\tau}}(\mathcal{D}_j)$ when the estimator is constructed based only on \mathcal{D}_j . Finally, we write $\ell_{n_k}^{(j)}(\boldsymbol{\beta}) = \sum_{i \in \mathcal{Q}_j} \ell_i(\boldsymbol{\beta})$ to denote the negative log-likelihood function based on \mathcal{D}_j .

While the results of this paper hold in a general likelihood based framework, for simplicity we state conditions at the population level for the generalized linear model (GLM) with canonical link. A much more general set of statements appear in the auxiliary lemmas upon which our main results are based. Under GLM with the canonical link, the response follows the distribution,

$$f_n(\mathbf{Y}; X, \boldsymbol{\beta}^*) = \prod_{i=1}^n f(Y_i; \eta_i^*) = \prod_{i=1}^n \left\{ c(Y_i) \exp \left[\frac{Y_i \eta_i^* - b(\eta_i^*)}{\phi} \right] \right\}, \quad (2.7)$$

where $\eta_i^* = \mathbf{X}_i^T \boldsymbol{\beta}^*$. The negative log-likelihood corresponding to (2.7) is given, up to an affine transformation, by

$$\ell_n(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n -Y_i \mathbf{X}_i^T \boldsymbol{\beta} + b(\mathbf{X}_i^T \boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n -Y_i \eta_i + b(\eta_i) = \frac{1}{n} \sum_{i=1}^n \ell_i(\boldsymbol{\beta}), \quad (2.8)$$

and the gradient and Hessian of $\ell_n(\boldsymbol{\beta})$ are respectively

$$\nabla \ell_n(\boldsymbol{\beta}) = -\frac{1}{n} \mathbf{X}^T (\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta})) \quad \text{and} \quad \nabla^2 \ell_n(\boldsymbol{\beta}) = \frac{1}{n} \mathbf{X}^T D(\mathbf{X}\boldsymbol{\beta}) \mathbf{X},$$

where $\boldsymbol{\mu}(\boldsymbol{\beta}) = (b'(\eta_1), \dots, b'(\eta_n))^T$ and $D(\boldsymbol{\beta}) = \text{diag}\{b''(\eta_1), \dots, b''(\eta_n)\}$. In this setting, $J(\boldsymbol{\beta}) = \mathbb{E}[b''(\mathbf{X}_1^T \boldsymbol{\beta}) \mathbf{X}_1 \mathbf{X}_1^T]$ and $J^* = \mathbb{E}[b''(\mathbf{X}_1^T \boldsymbol{\beta}^*) \mathbf{X}_1 \mathbf{X}_1^T]$.

3. Divide and Conquer Hypothesis Tests

In the context of the two classical testing frameworks, the Wald and Rao’s score tests, our objective is to construct a test statistic \bar{S}_n with low communication cost and a tractable limiting distribution F . From this statistic we define a test of size α of the null hypothesis, $H_0: \boldsymbol{\beta}_v^* = \boldsymbol{\beta}_v^H$, against the alternative, $H_1: \boldsymbol{\beta}_v^* \neq \boldsymbol{\beta}_v^H$, as a partition of the sample space described by

$$T_n^\alpha = \begin{cases} 0 & \text{if } |\bar{S}_n| \leq F^{-1}(1 - \alpha/2) \\ 1 & \text{if } |\bar{S}_n| > F^{-1}(1 - \alpha/2) \end{cases} \quad (3.1)$$

for a two sided test.

3.1. Two Divide and Conquer Wald Type Constructions

For the high dimensional linear model, Zhang and Zhang (2014), van de Geer et al. (2014) and Javanmard and Montanari (2014) propose methods for debiasing the Lasso estimator with a view to constructing high dimensional analogues of Wald statistics and confidence intervals for low-dimensional coordinates. As pointed out by Zhang and Zhang (2014), the debiased estimator does not impose the minimum signal condition used in establishing oracle properties of regularized estimators (Fan and Li, 2001; Fan and Lv, 2011; Loh and

Wainwright, 2015; Wang et al., 2014b; Zhang and Zhang, 2012) and hence has wider applicability than those inferences based on the oracle properties. The method of van de Geer et al. (2014) is appealing in that it accommodates a general penalized likelihood based framework, while the Javanmard and Montanari (2014) approach is appealing in that it optimizes asymptotic variance and requires a weaker condition than van de Geer et al. (2014) in the specific case of the linear model. We consider the DC analogues of Javanmard and Montanari (2014) and van de Geer et al. (2014) in Sections 3.1.1 and 3.1.2 respectively.

3.1.1. Lasso based Wald Test for the Linear Model—The linear model assumes

$$Y_i = \mathbf{X}_i^T \boldsymbol{\beta}^* + \varepsilon_i, \quad (3.2)$$

where $\{\varepsilon_i\}_{i=1}^n$ are i.i.d. with $\mathbb{E}(\varepsilon_i) = 0$ and variance σ^2 . For concreteness, we focus on a Lasso based method, but our procedure is also valid when other pilot estimators are used. We describe a modification of the bias correction method introduced in Javanmard and Montanari (2014) as a means to testing hypotheses on low dimensional coordinates of $\boldsymbol{\beta}^*$ via pivotal test statistics.

On each subset \mathcal{D}_j , we compute the debiased estimator of $\boldsymbol{\beta}^*$ as in Javanmard and Montanari (2014) as

$$\hat{\boldsymbol{\beta}}^d(\mathcal{D}_j) = \hat{\boldsymbol{\beta}}_{\text{Lasso}}^\lambda(\mathcal{D}_j) + \frac{1}{n_k} M^{(j)} (X^{(j)})^T (Y^{(j)} - X^{(j)} \hat{\boldsymbol{\beta}}_{\text{Lasso}}^\lambda(\mathcal{D}_j)), \quad (3.3)$$

where the superscript d is used to indicate the debiased version of the estimator,

$M^{(j)} = (\mathbf{m}_1^{(j)}, \dots, \mathbf{m}_d^{(j)})^T$ and \mathbf{m}_v is the solution of

$$\mathbf{m}_v^{(j)} = \underset{\mathbf{m}}{\operatorname{argmin}} \mathbf{m}^T \widehat{\boldsymbol{\Sigma}}^{(j)} \mathbf{m} \quad \text{s.t.} \quad \|\widehat{\boldsymbol{\Sigma}}^{(j)} \mathbf{m} - \mathbf{e}_v\|_\infty \leq \vartheta_1, \quad \|X^{(j)} \mathbf{m}\|_\infty \leq \vartheta_2. \quad (3.4)$$

The choice of tuning parameters ϑ_1 and ϑ_2 is discussed in Javanmard and Montanari (2014) and Zhao et al. (2014a) and they suggest to choose $\vartheta_1 \asymp \sqrt{\log d/n}$, $\vartheta_2 n^{-1/2} = \alpha(1)$. In the context of our DC procedure, ϑ_1 and ϑ_2 rely on k and should be chosen as $\vartheta_1 \asymp \sqrt{k \log d/n}$, $\vartheta_2 n^{-1/2} = \alpha(1)$, as quantified in Theorem 3.3. Above, $\widehat{\boldsymbol{\Sigma}}^{(j)} = n_k^{-1} \sum_{i \in \mathcal{D}_j} X_i^{(j)} X_i^{(j)T}$ is the sample covariance based on \mathcal{D}_j , whose population counterpart is $\boldsymbol{\Sigma} = \mathbb{E}(X_1 X_1^T)$ and $M^{(j)}$ is its regularized inverse. The second term in (3.3) is a bias correction term, while $\sigma^2 \mathbf{m}_v^{(j)T} \widehat{\boldsymbol{\Sigma}}^{(j)} \mathbf{m}_v^{(j)} / n_k$ is shown in Javanmard and Montanari (2014) to be the variance of the v^{th}

component of $\hat{\beta}^d(\mathcal{D}_j)$. The parameter ϑ_1 , which tends to zero, controls the bias of the debiased estimator (3.3) and the optimization in (3.4) minimizes the variance of the resulting estimator.

Solving d optimization problems in (3.4) increases an order of magnitude of computation complexity even for $k = 1$. Thus, it is necessary to appeal to the divide and conquer strategy to reduce the computation burden. This gives rise to the question how large k can be in order to maintain the same statistical properties as the whole sample one ($k = 1$).

Because our DC procedure gives rise to smaller samples, $\hat{\Sigma}$ is singular. This singularity does not pose a statistical problem but it does make the optimization problem ill-posed. To overcome the singularity in $\hat{\Sigma}$ and the resulting instability of the algorithm, we propose a change of variables. More specifically, noting that $M^{(j)}$ is not required explicitly, but rather the product $M^{(j)}(X^{(j)})^T$, we propose

$$\mathbf{b}_v^{(j)} = \underset{\mathbf{b}}{\operatorname{argmin}} \frac{\mathbf{b}^{(j)T} \mathbf{b}^{(j)}}{n_k} \quad \text{s.t.} \quad \left\| \frac{X^{(j)T} \mathbf{b}^{(j)}}{n_k} - \mathbf{e}_v \right\|_{\infty} \leq \vartheta_1, \quad \|\mathbf{b}^{(j)}\|_{\infty} \leq \vartheta_2, \quad (3.5)$$

from which we construct $M^{(j)}(X^{(j)})^T = B^T$, where $B = (\mathbf{b}_1, \dots, \mathbf{b}_d)$. The algorithm in equation (3.5) is crucial to the success of our procedure in practice.

The following conditions on the data generating process and the tail behavior of the design vectors are imposed in Javanmard and Montanari (2014). Both conditions are used to derive the theoretical properties of the DC Wald test statistic based on the aggregated debiased estimator, $\bar{\beta}^d = k^{-1} \sum_{j=1}^k \hat{\beta}^d(\mathcal{D}_j)$.

Condition 3.1: $\{(Y_i, X_i)\}_{i=1}^n$ are i.i.d. and Σ satisfies $0 < C_{\min} \lambda_{\min}(\Sigma) \lambda_{\max}(\Sigma) C_{\max}$.

Condition 3.2: The rows of X are sub-Gaussian with $\|X_i\|_{\psi_2} \leq \kappa, i = 1, \dots, n$.

Note that under the two conditions above, there exists a constant $\kappa_1 > 0$ such that

$$\|X_1 \Sigma^{-\frac{1}{2}}\|_{\psi_2} \leq \kappa_1. \text{ Without loss of generality, we set } \kappa_1 = \kappa. \text{ Our first main theorem provides}$$

the relative scaling of the various tuning parameters involved in the construction of $\bar{\beta}^d$.

Theorem 3.3: Suppose Conditions 2.1, 3.1 and 3.2 are fulfilled. Suppose $\mathbb{E}[\varepsilon_1^4] < \infty$ and choose ϑ_1, ϑ_2 and k such that $\vartheta_1 \asymp \sqrt{k \log d/n}, \vartheta_2 n^{-1/2} = \alpha(1)$ and $k = o((s \log d)^{-1} \sqrt{n})$. For any $v \in \{1, \dots, d\}$, we have

$$\sqrt{n} \frac{1}{k} \sum_{j=1}^k \frac{\hat{\beta}_v^d(\mathcal{D}_j) - \beta_v^*}{\hat{Q}_v^{(j)}} \rightsquigarrow N(0, \sigma^2), \quad (3.6)$$

where $\hat{Q}_v^{(j)} = (\mathbf{m}_v^{(j)T} \widehat{\Sigma}^{(j)} \mathbf{m}_v^{(j)})^{1/2}$.

Theorem 3.3 entertains the prospect of a divide and conquer Wald statistic of the form

$$\bar{S}_n = \sqrt{n} \frac{1}{k} \sum_{j=1}^k \frac{\hat{\beta}_v^d(\mathcal{D}_j) - \beta_v^H}{\bar{\sigma}(\mathbf{m}_v^{(j)T} \widehat{\Sigma}^{(j)} \mathbf{m}_v^{(j)})^{1/2}} \quad (3.7)$$

for β_v^* , where $\bar{\sigma}$ is an estimator for σ based on the k subsamples. On the left hand side of equation (3.7) we suppress the dependence on v to simplify notation. As an estimator for σ , a simple suggestion with the same computational complexity is $\bar{\sigma}$ where

$$\bar{\sigma}^2 = \frac{1}{k} \sum_{j=1}^k \hat{\sigma}^2(\mathcal{D}_j) \quad \text{and} \quad \hat{\sigma}^2(\mathcal{D}_j) = \frac{1}{n_k} \sum_{i \in \mathcal{I}_j} (Y_i^{(j)} - \mathbf{X}_i^{(j)T} \hat{\beta}_{\text{Lasso}}^\lambda(\mathcal{D}_j))^2. \quad (3.8)$$

One can use the refitted cross-validation procedure of Fan et al. (2012) to reduce the bias of the estimate. In Lemma 3.4 we show that with the scaling of k and λ required for the weak convergence results of Theorem 3.3, consistency of $\bar{\sigma}^2$ is also achieved.

Lemma 3.4: Suppose $\mathbb{E}[e_i | \mathbf{X}_i] = 0$ for all $i \in \{1, \dots, n\}$. Then with $\lambda \asymp \sqrt{k \log d/n}$ and $k = o(\sqrt{n}(s \log d)^{-1})$, $|\bar{\sigma}^2 - \sigma^2| = o_{\mathbb{P}}(1)$.

With Lemma 3.4 and Theorem 3.3 at hand, we establish in Corollary 3.5 the asymptotic distribution of \bar{S}_n under the null hypothesis $H_0: \beta_v^* = \beta_v^H$. This holds for each component $v \in \{1, \dots, d\}$.

Corollary 3.5: Suppose Conditions 3.1 and 3.2 are fulfilled, $\mathbb{E}[\varepsilon_1^4] < \infty$, and λ , ϑ_1 and ϑ_2 are chosen as $\lambda \asymp \sqrt{k \log d/n}$, $\vartheta_1 \asymp \sqrt{k \log d/n}$ and $\vartheta_2 n^{-1/2} = o(1)$. Then provided $k = o((s \log d)^{-1} \sqrt{n})$, under $H_0: \beta_v^* = \beta_v^H$, we have

$$\lim_{n \rightarrow \infty} |\mathbb{P}(\bar{S}_n \leq t) - \Phi(t)| = 0, \quad (3.9)$$

where $\Phi(\cdot)$ is the cdf of a standard normal distribution.

3.1.2. Wald Test in the Likelihood Based Framework—An alternative route to debiasing the Lasso estimator of β^* is the one proposed in van de Geer et al. (2014). Their so called desparsified estimator of β^* is more general than the debiased estimator of Javanmard and Montanari (2014) in that it accommodates generic estimators of the form (2.2) as pilot estimators, but the latter optimizes the variance of the resulting estimator. The desparsified estimator for subsample \mathcal{D}_j is

$$\hat{\beta}^d(\mathcal{D}_j) = \hat{\beta}^\lambda(\mathcal{D}_j) - \hat{\Theta}^{(j)} \nabla \ell_{n_k}^{(j)}(\hat{\beta}^\lambda(\mathcal{D}_j)), \quad (3.10)$$

where $\hat{\Theta}^{(j)}$ is a regularized inverse of the Hessian matrix of second order derivatives of $\ell_{n_k}^{(j)}(\beta)$ at $\hat{\beta}^\lambda(\mathcal{D}_j)$, denoted by $\hat{J}^{(j)} = \nabla^2 \ell_{n_k}^{(j)}(\hat{\beta}^\lambda(\mathcal{D}_j))$. We will make this explicit in due course.

The estimator resembles the classical one-step estimator (Bickel, 1975), but now in the high-dimensional setting via regularized inverse of the Hessian matrix $\hat{J}^{(j)}$, which reduces to the empirical covariance of the design matrix in the case of the linear model. From equation (3.10), the aggregated debiased estimator over the k subsamples is defined as

$$\bar{\beta}^d = k^{-1} \sum_{j=1}^k \hat{\beta}^d(\mathcal{D}_j).$$

We now use the nodewise Lasso (Meinshausen and Bühlmann, 2006) to approximately invert $\hat{J}^{(j)}$ via L_1 -regularization. The basic idea is to find the regularized invert row by row via a penalized L_1 -regression, which is the same as regressing the variable X_v on \mathbf{X}_{-v} but expressed in the sample covariance form. For each row $v \in 1, \dots, d$, consider the optimization

$$\hat{\kappa}_v(\mathcal{D}_j) = \underset{\kappa \in \mathbb{R}^{d-1}}{\operatorname{argmin}} \left(\hat{J}_{vv}^{(j)} - 2\hat{J}_{v,-v}^{(j)}\kappa + \kappa^T \hat{J}_{-v,-v}^{(j)}\kappa + 2\lambda_v \|\kappa\|_1 \right), \quad (3.11)$$

where $\hat{J}_{v,-v}^{(j)}$ denotes the v^{th} row of $\hat{J}^{(j)}$ without the $(v, v)^{\text{th}}$ diagonal element, and $\hat{J}_{-v,-v}^{(j)}$ is the principal submatrix without the v^{th} row and v^{th} column. Introduce

$$\hat{C}^{(j)} := \begin{pmatrix} 1 & -\hat{\kappa}_{1,2}(\mathcal{D}_j) & \dots & -\hat{\kappa}_{1,d}(\mathcal{D}_j) \\ -\hat{\kappa}_{2,1}(\mathcal{D}_j) & 1 & \dots & -\hat{\kappa}_{2,d}(\mathcal{D}_j) \\ \vdots & \vdots & \ddots & \vdots \\ -\hat{\kappa}_{d,1}(\mathcal{D}_j) & -\hat{\kappa}_{d,2}(\mathcal{D}_j) & \dots & 1 \end{pmatrix} \quad (3.12)$$

and $\widehat{\Xi}^{(j)} = \text{diag}(\widehat{\tau}_1(\mathcal{D}_j), \dots, \widehat{\tau}_d(\mathcal{D}_j))$, where $\widehat{\tau}_v(\mathcal{D}_j)^2 = \widehat{J}_{vv}^{(j)} - \widehat{J}_{v, -v}^{(j)} \widehat{\kappa}_v(\mathcal{D}_j)$. $\widehat{\Theta}^{(j)}$ in equation (3.10) is given by

$$\widehat{\Theta}^{(j)} = (\widehat{\Xi}^{(j)})^{-2} \widehat{C}^{(j)}, \quad (3.13)$$

and we define $\widehat{\Theta}_v^{(j)}$ as the transposed v^{th} row of $\widehat{\Theta}^{(j)}$.

Theorem 3.8 establishes the limit distribution of the term,

$$\bar{S}_n = \sqrt{n} \frac{1}{k} \sum_{j=1}^k \frac{\widehat{\beta}_v^d(\mathcal{D}_j) - \beta_v^H}{\sqrt{\widehat{\Theta}_{vv}^*}} \quad (3.14)$$

for any $v \in \{1, \dots, d\}$ under the null hypothesis $H_0: \beta_v^* = \beta_v^H$. This provides the basis for the statistical testing based on divide-and-conquer. We need the following condition. Recall that $J^* = E[\nabla_{\beta} \beta_n^l(\beta^*)]$ and consider the generalized linear model (2.7).

Condition 3.6: (i) $\{(Y_i, X_i)\}_{i=1}^n$ are i.i.d., $0 < C_{\min} \leq \lambda_{\min}(\Sigma) \leq \lambda_{\max}(\Sigma) \leq C_{\max}$, $\lambda_{\min}(J^*) > 0$, $\|J^*\|_{\max} < U_1 < \infty$. (ii) For some constant $M < \infty$, $\max_{1 \leq i \leq n} |X_i^T \beta^*| \leq M$ and $\max_{1 \leq i \leq n} \|X_i\|_{\infty} \leq M$. (iii) There exist finite constants $U_2, U_3 > 0$ such that $b''(\eta) < U_2$ and $b'''(\eta) < U_3$ for all $\eta \in \mathbb{R}$.

The same assumptions appear in van de Geer et al. (2014). In the case of the Gaussian GLM, the condition on $\lambda_{\min}(J^*)$ reduces to the requirement that the covariance of the design has minimal eigenvalue bounded away from zero, which is a standard assumption. We require $\|J^*\|_{\max} < \infty$ to control the estimation error of different functionals of J^* . The restriction in (ii) on the covariates and the projection of the covariates are imposed for technical simplicity; it can be extended to the case of exponential tails (see Fan and Song, 2010). Note that $\text{Var}(Y_i) = \phi b''(X_i^T \beta^*)$ where ϕ is the dispersion parameter in (2.7), so $b''(\eta) < U_2$ essentially implies an upper bound on the variance of the response. In fact, Lemma E.2 shows that $b''(\eta) < U_2$ can guarantee that the response is sub-Gaussian. $b'''(\eta) < U_3$ is used to derive the Lipschitz property of $b''(X_i^T \beta)$ with respect to β as shown in Lemma E.5. We emphasize that no requirement in Condition 3.6 is specific to the divide and conquer framework.

The assumption of bounded design in (ii) can be relaxed to the sub-Gaussian design. However, the price to pay is that the allowable number of subsets k is smaller than the bounded case, which means we need a larger sub-sample size. To be more precise, the order of maximum k for the sub-Gaussian design has an extra factor, which is a polynomial of $\sqrt{\log d}$, compared to the order for the bounded design. This logarithmic factor comes from

different Lipschitz properties of $b''(X_i^T \beta)$ in the two designs, which is fully explained in Lemma E.5 in the Supplementary Material. In the following theorems, we only present results for the case of bounded design for technical simplicity.

In addition, recalling that $\Theta^* = (J^*)^{-1}$, where $J^* := J(\beta^*) = \mathbb{E}[\nabla_{\beta\beta}^2 \ell_n(\beta^*)]$, we impose Condition 3.7 on Θ^* and its estimator $\hat{\Theta}$.

Condition 3.7: (i) $\min_{1 \leq v \leq d} \Theta_{vv}^* > \theta_{\min} > 0$. (ii) $\max_{1 \leq i \leq n} \|X_i^T \Theta^*\|_{\infty} \leq M$. (iii) For $v = 1, \dots, d$, whenever $\lambda_v \asymp \sqrt{k \log d/n}$ in (3.11), we have

$$\mathbb{P}\left(\|\hat{\Theta}_v - \Theta_v^*\|_1 \geq C s_1 \sqrt{\log d/n}\right) \leq d^{-1},$$

where C is a constant and s_1 is such that $\|\Theta_v^*\|_0 \lesssim s_1$ for all $v \in \{1, \dots, d\}$.

Part (i) of Corollary 3.7 ensures that the variances of each component of the debiased estimator exist, guaranteeing the existence of the Wald statistic. Parts (ii) and (iii) are imposed directly for technical simplicity. Results of this nature have been established under a similar set of assumptions in van de Geer et al. (2014) and Negahban et al. (2009) for convex penalties and in Wang et al. (2014a) and Loh and Wainwright (2015) for folded concave penalties.

As a step towards deriving the limit distribution of the proposed divide and conquer Wald statistic in the GLM framework, we establish the asymptotic behavior of the aggregated debiased estimator $\tilde{\beta}_v^d$ for every given $v \in [d]$.

Theorem 3.8: Under Conditions 2.1, 3.6 and 3.7, with $\lambda \asymp \sqrt{k \log d/n}$, we have

$$\tilde{\beta}_v^d - \beta_v^* = -\frac{1}{k} \sum_{j=1}^k \hat{\Theta}_v^{(j)T} \nabla \ell_{n_k}^{(j)}(\beta^*) + o_{\mathbb{P}}(n^{-1/2}) \quad (3.15)$$

for any $k \ll d$ satisfying $k = o\left(\left((s \vee s_1) \log d\right)^{-1} \sqrt{n}\right)$, where $\hat{\Theta}_v^{(j)}$ is the transposed v^{th} row of $\hat{\Theta}^{(j)}$.

The proof of Theorem 3.8 shows that for the Wald test procedure, the divide and conquer estimator $\tilde{\beta}_v^d$ is asymptotically as efficient as the full sample estimator $\hat{\beta}_v$, i.e.,

$$\lim_{n \rightarrow \infty} \frac{\text{Var}(\tilde{\beta}_v^d)}{\text{Var}(\hat{\beta}_v^d)} - 1 = 0.$$

A corollary of Theorem 3.8 provides the asymptotic distribution of the Wald statistic in equation (3.14) under the null hypothesis.

Corollary 3.9: Let \tilde{S}_n be as in equation (3.14), with Θ_{vv}^* replaced with an estimator $\tilde{\Theta}_{vv}$. Then under the conditions of Theorem 3.8 and $H_0: \beta_v^* = \beta_v^H$, provided $\|\tilde{\Theta}_{vv} - \Theta_{vv}^*\| = o_{\mathbb{P}}(1)$ under the scaling $k = o((s \vee s_1) \log d)^{-1} \sqrt{n}$, we have

$$\lim_{n \rightarrow \infty} \sup_{t \in \mathbb{R}} |\mathbb{P}(\tilde{S}_n \leq t) - \Phi(t)| = 0.$$

Remark 3.10: Although Theorem 3.8 and Corollary 3.9 are stated only for the GLM, their proofs are in fact an application of two more general results. Further details are available in Lemmas E.7 and E.8 in the Supplementary Material.

We return to the issue of estimating Θ_{vv}^* in Section 4, where we introduce a consistent estimator of Θ_{vv}^* that preserves the scaling of Theorem 3.8 and Corollary 3.9.

3.2. Divide and Conquer Score Test

In this section, we use $\nabla_v f(\beta)$ and $\nabla_{-v} f(\beta)$ to denote, respectively, the partial derivative of f with respect to β_v and the partial derivative vector of f with respect to β_{-v} . $\nabla_{vv}^2 f(\beta)$, $\nabla_{v,-v}^2 f(\beta)$, $\nabla_{-v,v}^2 f(\beta)$ and $\nabla_{-v,-v}^2 f(\beta)$ are analogously defined.

In the low dimensional setting (where d is fixed), Rao's score test of $H_0: \beta_v^* = \beta_v^H$ against $H_1: \beta_v^* \neq \beta_v^H$ is based on $\nabla_v \ell_n(\beta_v^H, \tilde{\beta}_{-v})$, where $\tilde{\beta}_{-v}$ is a constrained maximum likelihood estimator of β_{-v}^* , constructed as $\tilde{\beta}_{-v} = \operatorname{argmin}_{\beta_{-v}} \ell_n(\beta_v^H, \beta_{-v}) = \operatorname{argmax}_{\beta_{-v}} \{-\ell_n(\beta_v^H, \beta_{-v})\}$. If H_0 is false, imposing the constraint postulated by H_0 significantly violates the first order conditions from M-estimation with high probability; this is the principle underpinning the classical score test. Under regularity conditions, it can be shown (e.g. Cox and Hinkley, 1974) that

$$\sqrt{n}(\nabla_v \ell_n(\beta_v^H, \tilde{\beta}_{-v}))J_{v|-v}^* - 1/2 \rightsquigarrow N(0, 1),$$

where $J_{v|-v}^*$ is given by $J_{v|-v}^* = J_{v,v}^* - J_{v,-v}^* J_{-v,-v}^{*-1} J_{-v,v}^*$, with $J_{v,v}^*$, $J_{v,-v}^*$, $J_{-v,-v}^*$ and $J_{-v,v}^*$ the partitions of the information matrix $J^* = \mathcal{J}(\beta^*)$,

$$J(\boldsymbol{\beta}) = \begin{pmatrix} J_{v,v} & J_{v,-v} \\ J_{-v,v} & J_{-v,-v} \end{pmatrix} = \begin{pmatrix} \mathbb{E} \nabla_{v,v}^2 \ell_n(\boldsymbol{\beta}) & \mathbb{E} \nabla_{v,-v}^2 \ell_n(\boldsymbol{\beta}) \\ \mathbb{E} \nabla_{-v,v}^2 \ell_n(\boldsymbol{\beta}) & \mathbb{E} \nabla_{-v,-v}^2 \ell_n(\boldsymbol{\beta}) \end{pmatrix}. \quad (3.16)$$

The problems associated with the use of the classical score statistic in the presence of a high dimensional nuisance parameter are brought to light by Ning and Liu (2014), who propose a remedy via the decorrelated score. The problem stems from the inversion of the matrix $J_{-v,-v}^*$ in high dimensions. The decorrelated score is defined as

$$S(\boldsymbol{\beta}_v^*, \boldsymbol{\beta}_{-v}^*) = \nabla_v \ell_n(\boldsymbol{\beta}_v^*, \boldsymbol{\beta}_{-v}^*) - \mathbf{w}^{*T} \nabla_{-v} \ell_n(\boldsymbol{\beta}_v^*, \boldsymbol{\beta}_{-v}^*), \quad (3.17)$$

where $\mathbf{w}^{*T} = J_{v,-v}^* J_{-v,-v}^{*-1}$. For a regularized estimator $\hat{\mathbf{w}}$ of \mathbf{w}^* , to be defined below, we consider the score estimator

$$\hat{S}(\boldsymbol{\beta}_v^*, \hat{\boldsymbol{\beta}}_{-v}^\lambda) = \nabla_v \ell_n(\boldsymbol{\beta}_v^*, \hat{\boldsymbol{\beta}}_{-v}^\lambda) - \hat{\mathbf{w}}^T \nabla_{-v} \ell_n(\boldsymbol{\beta}_v^*, \hat{\boldsymbol{\beta}}_{-v}^\lambda). \quad (3.18)$$

Hence, provided \mathbf{w}^* is sufficiently sparse to avoid excessive noise accumulation, we are able to achieve consistency of $\hat{\mathbf{w}}$ under the high dimensional setting, ultimately giving rise to a tractable limit distribution of a suitable rescaling of $\hat{S}(\boldsymbol{\beta}_v^*, \hat{\boldsymbol{\beta}}_{-v}^\lambda)$. Since $\boldsymbol{\beta}_v^*$ is restricted under the null hypothesis, $H_0: \boldsymbol{\beta}_v^* = \boldsymbol{\beta}_v^H$, the statistic in (3.18) is accessible once H_0 is imposed. As Ning and Liu (2014) point out, \mathbf{w}^* is the solution to

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \mathbb{E} [\nabla_v \ell_n(\boldsymbol{\beta}_v^H, \boldsymbol{\beta}_{-v}^*) - \mathbf{w}^T \nabla_{-v} \ell_n(\boldsymbol{\beta}_v^H, \boldsymbol{\beta}_{-v}^*)]^2$$

under $H_0: \boldsymbol{\beta}_v^* = \boldsymbol{\beta}_v^H$.

Our divide and conquer score statistic under $H_0: \boldsymbol{\beta}_v^* = \boldsymbol{\beta}_v^H$ is

$$\bar{S}(\boldsymbol{\beta}_v^H) = \frac{1}{k} \sum_{j=1}^k \hat{S}^{(j)}(\boldsymbol{\beta}_v^H, \hat{\boldsymbol{\beta}}_{-v}^\lambda(\mathcal{D}_j)), \text{ where} \quad (3.19)$$

$\hat{S}^{(j)}(\boldsymbol{\beta}_v, \hat{\boldsymbol{\beta}}_{-v}^\lambda(\mathcal{D}_j)) = \nabla_v \ell_{n_k}^{(j)}(\boldsymbol{\beta}_v, \hat{\boldsymbol{\beta}}_{-v}^\lambda(\mathcal{D}_j)) - \hat{\mathbf{w}}(\mathcal{D}_j)^T \nabla_{-v} \ell_{n_k}^{(j)}(\boldsymbol{\beta}_v, \hat{\boldsymbol{\beta}}_{-v}^\lambda(\mathcal{D}_j))$ and we estimate \mathbf{w}^* using the Dantzig selector of Candes and Tao (2007)

$$\begin{aligned} \hat{\mathbf{w}}(\mathcal{D}_j) &= \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|_1, \\ \text{s.t. } &\left\| \nabla_{-v, v}^2 \ell_{n_k}^{(j)}(\hat{\beta}_v^\lambda(\mathcal{D}_j), \hat{\beta}_{-v}^\lambda(\mathcal{D}_j)) - \mathbf{w}^T \nabla_{-v, -v}^2 \ell_{n_k}^{(j)}(\hat{\beta}_v^\lambda(\mathcal{D}_j), \hat{\beta}_{-v}^\lambda(\mathcal{D}_j)) \right\|_\infty \leq \mu. \end{aligned}$$

Theorem 3.11: Let $\hat{\mathcal{J}}_{v|-v}$ be a consistent estimator of $J_{v|-v}^*$ and

$$S^{(j)}(\beta_v^H, \beta_{-v}^*) = \nabla_{v, v} \ell_{n_k}^{(j)}(\beta_v^H, \beta_{-v}^*) - \mathbf{w}^{*T} \nabla_{-v, -v} \ell_{n_k}^{(j)}(\beta_v^H, \beta_{-v}^*).$$

Suppose $\|\mathbf{w}^*\|_1 \lesssim s_1$ and Conditions 2.1 and 3.6 are fulfilled. Then under $H_0: \beta_v^* = \beta_v^H$ with $\lambda \asymp \mu \asymp \sqrt{k \log d/n}$,

$$\begin{aligned} \sqrt{n} \bar{S}(\beta_v^H) &= \sqrt{n} \frac{1}{k} \sum_{j=1}^k S^{(j)}(\beta_v^H, \beta_{-v}^*) + o_{\mathbb{P}}(1) \\ \text{and } \lim_{n \rightarrow \infty} \sup_{t \in \mathbb{R}} &\left| \mathbb{P}\left(\sqrt{n} \cdot \bar{S}(\beta_v^H) \hat{\mathcal{J}}_{v|-v}^{-1/2} \leq t\right) - \Phi(t) \right| = 0, \end{aligned}$$

for any $k \ll d$ satisfying $k = o\left(\left((s \vee s_1) \log d\right)^{-1} \sqrt{n}\right)$, where $\bar{S}(\beta_v^H)$ is defined in equation (3.19).

Remark 3.12: By the definition of \mathbf{w}^* and the block matrix inversion formula for $\Theta^* = (\mathcal{J}^*)^{-1}$, sparsity of \mathbf{w}^* is implied by sparsity of Θ^* as assumed in van de Geer et al. (2014) and Condition 3.7 of Section 3.1.2. In turn, $\|\mathbf{w}^*\|_0 \lesssim s_1$ implies $\|\mathbf{w}^*\|_1 \lesssim s_1$ provided that the elements of \mathbf{w}^* are bounded.

Remark 3.13: Although Theorem 3.11 is stated in the penalized GLM setting, the result holds more generally; further details are available in Lemma E.13 in the Supplementary Material.

To maintain the same computational complexity, an estimator of the conditional information needs to be constructed using a DC procedure. For this, we propose to use

$$\bar{J}_{v|-v} = k^{-1} \sum_{j=1}^k \left(\nabla_{v, v}^2 \ell_{n_k}^{(j)}(\bar{\beta}_v^d, \bar{\beta}_{-v}) - \bar{\mathbf{w}}^T \nabla_{-v, -v}^2 \ell_{n_k}^{(j)}(\bar{\beta}_v^d, \bar{\beta}_{-v}) \right),$$

where the divide and conquer estimator $\bar{\beta}_v^d = k^{-1} \sum_{j=1}^k \hat{\beta}_v^d(\mathcal{D}_j)$, $\bar{\beta}_{-v} = k^{-1} \sum_{j=1}^k \hat{\beta}_{-v}^\lambda(\mathcal{D}_j)$ and $\bar{\mathbf{w}} = k^{-1} \sum_{j=1}^k \hat{\mathbf{w}}(\mathcal{D}_j)$. Note that for certain v , the communication cost for calculating $\bar{J}_{v|-v}$ is not high, since all the involved quantities $\left\{ \nabla_{v, v}^2 \ell_{n_k}^{(j)}(\bar{\beta}_v^d, \bar{\beta}_{-v}) \right\}_{j=1}^k$, $\left\{ \nabla_{-v, -v}^2 \ell_{n_k}^{(j)}(\bar{\beta}_v^d, \bar{\beta}_{-v}) \right\}_{j=1}^k$

and $\{\hat{w}(\mathcal{D}_j)\}_{j=1}^k$ are scalars, d -dimensional vectors and d -dimensional vectors respectively.

The communication cost is thus of order $\mathcal{O}(kd)$. We do not communicate the entire huge hessian matrix here.

Lemma 3.14: Suppose $\|w^*\|_1 = \mathcal{O}(s_1)$ and Conditions 2.1 and 3.6 are fulfilled. Then for any $k \ll d$ satisfying $k = o\left(\left((s \vee s_1) \log d\right)^{-1} \sqrt{n}\right)$, $|\bar{J}_{v|-v} - J_{v|-v}^*| = o_{\mathbb{P}}(1)$.

By Lemma 3.14, we show that $\bar{J}_{v|-v}$ is consistent and can be applied to Theorem 3.11.

4. Accuracy of Distributed Estimation

This section focuses on high-dimensional ($d \gg n$) divide-and-conquer estimators for linear and generalized linear models. As explained below Theorem 3.8 in Section 3, the efficiency loss from the divide-and-conquer process is asymptotically zero. This motivates us to consider $\|\bar{\beta}^d - \hat{\beta}^d\|$, the loss incurred by the divide and conquer strategy in comparison with the practically unavailable full sample debiased estimator $\hat{\beta}^d$, where $\|\cdot\|$ is certain norm. Indeed, it turns out that, for k not too large, $\bar{\beta}^d - \hat{\beta}^d$ appears only as a higher order term in the decomposition of $\bar{\beta}^d - \beta^*$ and thus $\|\bar{\beta}^d - \hat{\beta}^d\|$ is negligible compared to the statistical error, $\|\hat{\beta}^d - \beta^*\|$. In other words, the divide-and-conquer errors are statistically negligible.

Compared with calculating the full sample debiased Lasso estimator, our proposed DC strategy enjoys computational advantages since it is highly parallel and each subsample problem has a much smaller scale than the full sample problem given a suitably large k . However, relative to just the full sample penalized M-estimator (e.g., Lasso), distributed point estimation does not entail a computational gain like distributed testing, since our distributed algorithm requires debiasing each component of the Lasso estimator and hence brings high expense of computation. The bottleneck of computation of our DC procedure comes from the d extra debiasing steps. To mitigate this problem, we can actually debias each component of $\hat{\beta}$ in a parallel fashion. According to the optimization procedures (3.4) and (3.11), debiasing one component of the Lasso estimator is entirely independent of the debiasing of another component. Therefore, as long as each branch computer in the cluster shares the sub-dataset \mathcal{D}_j and the Lasso estimator $\hat{\beta}^{(j)}$, they can work in parallel and collectively return to a central server all the components of the debiased Lasso estimator. This parallelization reduces the time complexity significantly.

When the minimum signal strength is sufficiently strong, thresholding $\bar{\beta}^d$ achieves exact support recovery, motivating a refitting procedure based on the low dimensional selected variables. As a means to understanding the theoretical properties of this refitting procedure, as well as for independent interest, we develop new theories and methodologies for the low dimensional ($d < n$) linear and generalized linear models in Appendixes A and B in the Supplementary Material respectively. We show that simple averaging of low dimensional OLS or GLM estimators (denoted uniformly as $\hat{\beta}^{(j)}$, without superscript d as debiasing is not necessary) suffices to preserve the statistical error, i.e., achieving the same statistical accuracy as the estimator based on the full sample. This is because, in contrast to the high dimensional setting, parameters are not penalized in the low dimensional case. With $\bar{\beta}$ the

average of $\hat{\beta}^{(j)}$ over the k machines and $\hat{\beta}$ the full sample counterpart ($k = 1$), we derive the rate of convergence of $\|\bar{\beta} - \hat{\beta}\|_2$. Refitted estimation using only the selected covariates allows us to eliminate the $\log d$ term in the statistical rate of convergence of the estimator under high-dimensional settings. We present theoretical results on the refitting estimation as corollaries to the low-dimensional regression results in Appendixes A and B in the Supplementary Material.

4.1. The High-Dimensional Linear Model

Recall that the high dimensional DC estimator is $\bar{\beta}^d = k^{-1} \sum_{j=1}^k \hat{\beta}^d(\mathcal{D}_j)$, where $\hat{\beta}^d(\mathcal{D}_j)$ for $1 \leq j \leq k$ is the debiased estimator defined in (3.3). We also denote the debiased Lasso estimator using the entire dataset as $\hat{\beta}^d = \hat{\beta}^d(\cup_{j=1}^k \mathcal{D}_j)$. The following lemma shows that not only is $\bar{\beta}^d$ asymptotically normal, it approximates the full sample estimator $\hat{\beta}^d$ so well that it has the same statistical error as $\hat{\beta}^d$ provided the number of subsamples k is not too large.

Lemma 4.1: Consider the linear model (3.2). Under the Conditions 3.1 and 3.2, if λ , ϑ_1 and ϑ_2 are chosen as $\lambda \asymp \sqrt{k \log d/n}$, $\vartheta_1 \asymp \sqrt{k \log d/n}$ and $\vartheta_2 n^{-1/2} = \alpha(1)$, we have with probability $1 - c/d$,

$$\|\bar{\beta}^d - \hat{\beta}^d\|_\infty \leq C \frac{sk \log d}{n} \text{ and } \|\bar{\beta}^d - \beta^*\|_\infty \leq C \left(\sqrt{\frac{\log d}{n}} + \frac{sk \log d}{n} \right). \quad (4.1)$$

Remark 4.2: The term $\sqrt{\log d/n}$ in (4.1) is the estimation error of $\|\hat{\beta}^d - \beta^*\|_\infty$, while the term $(sk \log d)/n$ is the rate of the distance between the divide and conquer estimator and the full sample estimator. Lemma 4.1 does not rely on any specific choice of k . However, in order for the aggregated estimator $\bar{\beta}^d$ to attain the same $\|\cdot\|_\infty$ norm estimation error as the full sample Lasso estimator, $\hat{\beta}_{\text{Lasso}}$, the required scaling is $k = O(\sqrt{n/(s^2 \log d)})$. This is a weaker scaling requirement than that of Theorem 3.3 because the latter entails a guarantee of asymptotic normality, which is a stronger result. It is for the same reason that our estimation results only require $O(\cdot)$ scaling whilst those for testing require $\alpha(\cdot)$ scaling.

Rosenblatt and Nadler (2016) show that in the high-dimensional regime where $d/n_k \rightarrow \kappa \in (0, 1)$, the divide and conquer procedure suffers from first-order accuracy loss. This seems a contradiction to our result, since our dimension is even higher than their context, but we have no first-order accuracy loss while averaging debiased estimators based on subsamples, as long as we have an appropriate number of data splits. In fact, in the highdimensional sparse linear regression, the intrinsic dimension is the sparsity s rather than d , which is regarded instead as the ambient dimension. The sparsity assumption changes the original high-dimensional problem to be an intrinsically low-dimensional one and thus allows us to escape from any first-order accuracy loss of the divide and conquer procedure. Given $s = \alpha(n_k)$, we can treat high-dimensional sparse linear regression approximately as the classical linear regression setting where $d = \alpha(n_k)$. Hence we expect no first-order accuracy loss from the divide and conquer procedure here.

Although $\tilde{\beta}^d$ achieves the same rate as the Lasso estimator under the infinity norm, it cannot achieve the minimax rate in ℓ_2 norm since it is not a sparse estimator. To obtain an estimator with the ℓ_2 minimax rate, we sparsify $\tilde{\beta}^d$ by hard thresholding. For any $\beta \in \mathbb{R}^d$, define the hard thresholding operator \mathcal{T}_ν such that the j^{th} entry of $\mathcal{T}_\nu(\beta)$ is

$$[\mathcal{T}_\nu(\beta)]_j = \beta_j \mathbb{1}\{|\beta_j| \geq \nu\}, \text{ for } 1 \leq j \leq d. \quad (4.2)$$

According to (4.1), if $\beta_j^* = 0$, we have $|\tilde{\beta}_j^d| \leq C(\sqrt{\log d/n} + sk \log d/n)$ with high probability. The following theorem characterizes the estimation error, $\|\mathcal{T}_\nu(\tilde{\beta}^d) - \beta^*\|_2$, and divide and conquer error, $\|\mathcal{T}_\nu(\tilde{\beta}^d) - \mathcal{T}_\nu(\hat{\beta}^d)\|_2$, of the thresholded estimator $\mathcal{T}_\nu(\tilde{\beta}^d)$.

Theorem 4.3: Under the linear model (3.2), suppose Conditions 3.1 and 3.2 are fulfilled and choose $\lambda \asymp \sqrt{k \log d/n}$, $\vartheta_1 \asymp \sqrt{k \log d/n}$ and $\vartheta_2 n^{-1/2} = \alpha(1)$. Take the parameter of the hard threshold operator in (4.2) as $\nu = C_0 \sqrt{\log d/n}$ for some sufficiently large constant C_0 . If the number of subsamples satisfies $k = O(\sqrt{nl}(s^2 \log d))$, for large enough d and n , we have with probability $1 - c/d$,

$$\begin{aligned} \|\mathcal{T}_\nu(\tilde{\beta}^d) - \mathcal{T}_\nu(\hat{\beta}^d)\|_2 &\leq C \frac{s^{3/2} k \log d}{n}, \quad \|\mathcal{T}_\nu(\tilde{\beta}^d) - \beta^*\|_\infty \leq C \sqrt{\frac{\log d}{n}} \\ \text{and } \|\mathcal{T}_\nu(\tilde{\beta}^d) - \beta^*\|_2 &\leq C \sqrt{\frac{s \log d}{n}}. \end{aligned}$$

Remark 4.4: In fact, in the proof of Theorem 4.3, we show that if the thresholding parameter ν satisfies $\nu \geq \|\tilde{\beta}^d - \beta^*\|_\infty$, we have $\|\mathcal{T}_\nu(\tilde{\beta}^d) - \beta^*\|_2 \leq 2\sqrt{2s} \cdot \nu$; it is for this reason that we choose $\nu \asymp \sqrt{\log d/n}$. Unfortunately, the constant is difficult to choose in practice. In the following paragraphs we propose a practical method to select the tuning parameter ν .

Let $(M^{(j)} X^{(j)T})_{\ell}$ denote the transposed ℓ^{th} row of $M^{(j)} X^{(j)T}$. Inspection of the proof of Theorem 3.3 reveals that the leading term of $\sqrt{n} \|\tilde{\beta}^d - \beta^*\|_\infty$ satisfies

$$T_0 = \max_{1 \leq \ell \leq d} \frac{1}{\sqrt{k}} \sum_{j=1}^k \frac{1}{\sqrt{n_k}} (M^{(j)} X^{(j)T})_{\ell}^T e^{(j)}.$$

Chernozhukov et al. (2013) propose the Gaussian multiplier bootstrap to estimate the quantile of T_0 . Let $\{\xi_i\}_{i=1}^n$ be i.i.d. standard normal random variable independent of $\{(Y_i, X_i)\}_{i=1}^n$. Consider the statistic

$$W_0 = \max_{1 \leq \ell \leq d} \frac{1}{\sqrt{k}} \sum_{j=1}^k \frac{1}{\sqrt{n_k}} (M^{(j)} X^{(j)T})^T \ell(\hat{\boldsymbol{\epsilon}}^{(j)} \circ \boldsymbol{\xi}^{(j)}),$$

where $\hat{\boldsymbol{\epsilon}}^{(j)} \in \mathbb{R}^{nk}$ is an estimator of $\boldsymbol{\epsilon}^{(j)}$ such that for any $i \in \mathcal{Q}_j$, $\hat{\epsilon}_i^{(j)} = Y_i^{(j)} - X_i^{(j)} \hat{\boldsymbol{\beta}}^{(j)}$, and $\boldsymbol{\xi}^{(j)}$ is a subvector of $\{\xi_i\}_{i=1}^n$ with indices in \mathcal{Q}_j . Recall that “ \circ ” denotes the Hadamard product. The α -quantile of W_0 conditioning on $\{Y_i, X_i\}_{i=1}^n$ is defined as $c_{W_0}(\alpha) = \inf\{t / \mathbb{P}(W_0 \leq t | Y, X) = \alpha\}$. We estimate $c_{W_0}(\alpha)$ by Monte-Carlo and thus choose $\nu_0 = c_{W_0}(\alpha) / \sqrt{n}$. This choice ensures

$$\left\| \mathcal{T}_{\nu_0}(\tilde{\boldsymbol{\beta}}^d) - \boldsymbol{\beta}^* \right\|_2 = O_{\mathbb{P}}(\sqrt{s \log d/n}),$$

which coincides with the ℓ_2 convergence rate of the Lasso.

Remark 4.5: Lemma 4.1 and Theorem 4.3 show that if the number of subsamples satisfies $k = o(\sqrt{n/(s^2 \log d)})$, $\|\tilde{\boldsymbol{\beta}}^d - \hat{\boldsymbol{\beta}}^d\|_\infty = o_{\mathbb{P}}(\sqrt{\log d/n})$ and $\|\mathcal{T}_{\nu}(\tilde{\boldsymbol{\beta}}^d) - \mathcal{T}_{\nu}(\hat{\boldsymbol{\beta}}^d)\|_2 = o_{\mathbb{P}}(\sqrt{s \log d/n})$, and thus the error incurred by the divide and conquer procedure is negligible compared to the statistical minimax rate. The reason for this contraction phenomenon is that $\tilde{\boldsymbol{\beta}}^d$ and $\hat{\boldsymbol{\beta}}^d$ share the same leading term in their Taylor expansions around $\boldsymbol{\beta}^*$. The difference between them is only the difference of two remainder terms which has a smaller order than the leading term. We uncover a similar phenomenon in the low dimensional case covered in Appendix A in the Supplementary Material. However, in the low dimensional case ℓ_2 norm consistency is automatic while the high dimensional case requires an additional thresholding step to guarantee sparsity and, consequently, ℓ_2 norm consistency.

4.2. The High-Dimensional Generalized Linear Model

We generalize the DC estimation of the linear model to GLM. Recall that $\hat{\boldsymbol{\beta}}^d(\mathcal{D}_j)$ is the de-biased estimator defined in (3.10) and the aggregated estimator is $\tilde{\boldsymbol{\beta}}^d = k^{-1} \sum_{j=1}^k \hat{\boldsymbol{\beta}}^d(\mathcal{D}_j)$. We still denote $\hat{\boldsymbol{\beta}}^d = \hat{\boldsymbol{\beta}}^d(\cup_{j=1}^k \mathcal{D}_j)$. The next lemma bounds the error incurred by splitting the sample and the statistical rate of convergence of $\tilde{\boldsymbol{\beta}}^d$ in terms of the infinity norm.

Lemma 4.6: Consider the generalized linear model (2.7) with canonical link. Under Conditions 2.1, 3.6 and 3.7, for $\hat{\boldsymbol{\beta}}^\lambda$ with $\lambda \asymp \sqrt{k \log d/n}$, we have with probability $1 - c/d$, there exists a constant $C > 0$, such that

$$\|\tilde{\boldsymbol{\beta}}^d - \hat{\boldsymbol{\beta}}^d\|_\infty \leq C \frac{(s \vee s_1)^k \log d}{n}, \|\tilde{\boldsymbol{\beta}}^d - \boldsymbol{\beta}^*\|_\infty \leq C \left(\sqrt{\frac{\log d}{n}} + \frac{(s \vee s_1)^k \log d}{n} \right).$$

Remark 4.7: The term $\sqrt{\log d/n}$ in above is the estimation error of $\|\hat{\beta}^d - \beta^*\|_\infty$, while the error term $(s \vee s_1)k \log d/n$ is attributable to the distance between $\bar{\beta}^d$ and $\hat{\beta}^d$.

Applying a similar thresholding step as in the linear model, we quantify the ℓ_2 -norm estimation error, $\|\mathcal{T}_\nu(\hat{\beta}^d) - \beta^*\|_2$ and the distance between the divide and conquer estimator and full sample estimator $\|\mathcal{T}_\nu(\bar{\beta}^d) - \mathcal{T}_\nu(\hat{\beta}^d)\|_2$.

Theorem 4.8: For the GLM (2.7), under Conditions 2.1 – 3.7, choose $\lambda \asymp \sqrt{k \log d/n}$ and $\lambda_\nu \asymp \sqrt{k \log d/n}$. Take the parameter of the hard threshold operator in (4.2) as $\nu = C_0 \sqrt{\log d/n}$ for some sufficiently large constant C_0 . If the number of subsamples satisfies $k = O(\sqrt{n/(s \vee s_1)^2 \log d})$, for large enough d and n , we have with probability $1 - c/d$,

$$\|\mathcal{T}_\nu(\bar{\beta}^d) - \mathcal{T}_\nu(\hat{\beta}^d)\|_2 \leq C \frac{(s \vee s_1) s^{1/2} k \log d}{n}, \quad \|\mathcal{T}_\nu(\bar{\beta}^d) - \beta^*\|_\infty \leq C \sqrt{\frac{\log d}{n}} \quad (4.3)$$

$$\text{and } \|\mathcal{T}_\nu(\bar{\beta}^d) - \beta^*\|_2 \leq C \sqrt{s \log d/n}.$$

Remark 4.9: As in the case of the linear model, Theorem 4.8 reveals that the loss incurred by the divide and conquer procedure is negligible compared to the statistical minimax estimation error provided $k = o(\sqrt{n/(s_1 \vee s)^2 s \log d})$.

A similar proof strategy to that of Theorem 4.8 allows us to construct an estimator of $\Theta_{\nu\nu}^*$ that achieves the required consistency with the scaling of Corollary 3.9. Our estimator is $\tilde{\Theta}_{\nu\nu} = [\mathcal{T}_\zeta(\bar{\Theta})]_{\nu\nu}$, where $\bar{\Theta} = k^{-1} \sum_{j=1}^k \hat{\Theta}^{(j)}$ and $\mathcal{T}_\zeta(\cdot)$ is the thresholding operator defined in equation (4.2) with $\zeta = C_1 \sqrt{\log d/n}$ for some sufficiently large constant C_1 .

Corollary 4.10: Under the conditions and scaling of Theorem 3.8, $|\tilde{\Theta}_{\nu\nu} - \Theta_{\nu\nu}^*| = o_{\mathbb{P}}(1)$.

Substituting this estimator in Corollary 3.9 delivers a practically implementable test statistic based on $k = o((s \vee s_1) \log d)^{-1} \sqrt{n}$ subsamples.

Remark 4.11: Notice that point estimation requires less stringent scaling of k than hypothesis testing in both the linear and generalized linear models. This is because the testing and estimation require different rates for the higher order term in the decomposition

$$\sqrt{n}(\bar{\beta}^d - \beta^*) = \mathbf{Z} + \mathbf{\Delta},$$

where \mathbf{Z} is the leading term contributing to the asymptotic normality of $\sqrt{n}(\bar{\beta}^d - \beta^*)$. For hypothesis testing, we need $\|\mathbf{\Delta}/\sqrt{n}\|_\infty = o_{\mathbb{P}}(1/\sqrt{n})$ to guarantee the asymptotic normality. For

estimation, we need $\|\Delta\sqrt{n}\|_\infty = o_p(\sqrt{\log d/n})$ to match the minimax rate of $\|\tilde{\beta}^d - \beta^*\|_\infty$.

Therefore, the number of splits k for testing is more stringent by a factor of $1/\sqrt{\log d}$ than in estimation.

5. Simulations

In this section, we illustrate and validate our theoretical findings through simulations. For hypothesis testing, we use QQ plots to compare the distribution of p -values for divide and conquer test statistics to their theoretical uniform distribution. We also investigate the estimated type I error and power of the divide and conquer tests. For estimation, we validate our claim in the beginning of Section 4 that the loss incurred by the divide and conquer strategy is negligible compared with the statistical error of the corresponding full sample estimator in the high dimensional case. Specifically, we compare the performance of the divide and conquer thresholding estimator of Section 4.1 with the full sample Lasso and the average Lasso over subsamples. An analogous empirical verification of the theory is performed for the low dimensional case as well; we put it in Appendixes C and D of the Supplementary Material.

5.1. Results on Hypothesis Testing

We explore the probability of rejection of a null hypothesis of the form $H_0: \beta_1^* = 0$ when data $(Y_i, X_i)_{i=1}^n$ are generated according to the linear model,

$$Y_i = X_i^T \beta^* + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma_\varepsilon^2),$$

where $\sigma_\varepsilon^2 = 1$ and

$$\beta^* = (\beta_1^*, \underbrace{0, \dots, 0}_{d-s-1}, \underbrace{1, \dots, 1}_s)^T,$$

where $d = 5000$ and $s = 3$. In each Monte Carlo replication, we split the initial sample of size n into k subsamples of size n/k . In particular we choose $n = 5000$ and $k \in \{1, 2, 5, 10, 20, 25, 40, 50, 100, 200, 500\}$. The number of Monte Carlo replications is 500. Using $\hat{\beta}_{\text{Lasso}}$ as a preliminary estimator of β^* , we construct Wald and Rao's score test statistics as described in Sections 3.1.2 and 3.2 respectively.

Panels (A) and (B) of Figure 1 are QQ plots of the p -values of the divide and conquer Wald and score test statistics under the null hypothesis against the theoretical quantiles of the uniform $[0,1]$ distribution for eight different values of k . For both test constructions, the distributions of the p -values are close to uniform and remain so as we split the data set. When $k = 100$, the distribution of the corresponding p -values deviates from the uniform distribution visibly, as expected from the theory developed in Sections 3.1.2 and 3.2. Panel

(A) of Figure 2 shows that, for both test constructions, when the number of splits $k = 50$, the empirical level of the test is close to both the nominal $\alpha = 0.05$ level and the level of the full sample oracle OLS estimator which knows the true support of β^* . On the other hand, the type I error increases dramatically when k is larger than 50. This is consistent with asymptotic normality of the test statistics we established when k is controlled appropriately. Panel (B) of Figure 2 displays the power of the test for two different signal strengths, $\beta_1^* = 0.05$ and 0.06 . We see that the power for the Score and Wald tests improves when the signal strength goes from 0.05 to 0.06 . In addition, we find that the power is high regardless of how large k is. However, Figure 2(A) shows that the Type I error is large when k is large, which makes the tests invalid. Therefore, these results illustrate that the Type I and II errors are controllable when the number of splits k is relatively small. We also record the wall time for computation for these k 's in Table 1. The wall time is computed by taking the maximal time taken for each split and averaged over replications.

5.2. Results on Estimation

In this section, we turn our attention to experimental validation of our divide and conquer estimation theory, focusing first on the low dimensional case and then on the high dimensional case.

5.2.1. The High Dimensional Linear Model—We now consider the same setting of Section 5.1 with $n = 5000$, $d = 5000$ and $\beta_j^* = 10$ for all j in the support of β^* . In this context, we analyze the performance of the thresholded averaged debiased estimator of Section 4.1. Figure 3(A) depicts the average over 100 Monte Carlo replications of $\|b - \beta^*\|_2$ for three different estimators: debiased divide-and-conquer $b = \mathcal{T}_\nu(\bar{\beta}^d)$, the Lasso estimator based on the whole sample $b = \hat{\beta}_{\text{Lasso}}$ and the estimator obtained by naïvely averaging the Lasso estimators from the k subsamples $b = \bar{\beta}_{\text{Lasso}}$. The parameter ν is taken as $\nu = \sqrt{\log d/n}$ in the specification of $\mathcal{T}_\nu(\bar{\beta}^d)$. As expected, the performance of $\bar{\beta}_{\text{Lasso}}$ deteriorates sharply as k increases. $\mathcal{T}_\nu(\bar{\beta}^d)$ outperforms $\hat{\beta}_{\text{Lasso}}$ as long as k is not too large. This is expected because, for sufficiently large signal strength, both $\hat{\beta}_{\text{Lasso}}$ and $\mathcal{T}_\nu(\bar{\beta}^d)$ recover the correct support, however $\mathcal{T}_\nu(\bar{\beta}^d)$ is unbiased for those β_j^* in the support of β^* , whilst $\hat{\beta}_{\text{Lasso}}$ is biased. Figure 3(B) shows the error incurred by the divide and conquer procedure $\|\mathcal{T}_\nu(\bar{\beta}^d) - \mathcal{T}_\nu(\hat{\beta}^d)\|_2$ relative to the statistical error of the full sample estimator, $\|\mathcal{T}_\nu(\bar{\beta}^d) - \beta^*\|_2$, for four different scalings of k . We observe that, with $k \asymp O(\sqrt{nl/s^2 \log d})$, the relative error incurred by the divide and conquer procedure can hardly converge. This is consistent with Theorem 4.3. Given the lower bound of statistical error of the full sample Lasso estimator $\hat{\beta}$. From Theorem 4.3 we derive that

$$\frac{E\|\mathcal{T}_\nu(\bar{\beta}^d) - \mathcal{T}_\nu(\hat{\beta}^d)\|_2^2}{E\|\mathcal{T}_\nu(\bar{\beta}^d) - \beta^*\|_2^2} \leq \frac{s^2 k^2 \log d}{n}.$$

When $k \asymp O(\sqrt{n/s^2 \log d})$, the righthand side is an $O(1)$ term. Therefore the line with inverted triangles in Figure 3(B) implies that the statistical error rate we developed in Theorem 4.3 is tight. We also record the wall time for estimation computation for these k 's in Table 1. The wall time is computed by taking the maximal time taken for each splits and averaged over replications. We notice that the computation time decreases with k at first due to the parallel algorithm. However, for the score test and split Lasso, the time becomes increasing when k is large, this is because the computation time to aggregate results from different splits is no longer negligible for very large k 's.

6. Discussion

With the advent of the data revolution comes the need to modernize the classical statistical tool kit. For very large scale datasets, distribution of data across multiple machines is the only practical way to overcome storage and computational limitations. It is thus essential to build aggregation procedures for conducting inference based on the combined output of multiple machines. We successfully achieve this objective, deriving divide and conquer analogues of the Wald and score statistics and providing statistical guarantees on their performance as the number of sample splits grows to infinity with the full sample size. Tractable limit distributions of each DC test statistic are derived. These distributions are valid as long as the number of subsamples, k , does not grow too quickly. In particular, $k = o((s \vee s_1) \log d)^{-1} \sqrt{n}$ is required in a general likelihood based framework. If k grows faster than $((s \vee s_1) \log d)^{-1} \sqrt{n}$, remainder terms become nonnegligible and contaminate the tractable limit distribution of the leading term. When attention is restricted to the linear model, a faster growth rate of $k = o(s \log d)^{-1} \sqrt{n}$ is allowed.

The divide and conquer strategy is also successfully applied to estimation of regression parameters. We obtain the rate of the loss incurred by the divide and conquer strategy. Based on this result, we derive an upper bound on the number of subsamples for preserving the statistical error. For low-dimensional models, simple averaging is shown to be effective in preserving the statistical error, so long as $k = O(n/d)$ for the linear model and $k = O(\sqrt{n}/d)$ for the generalized linear model. For high-dimensional models, the debiased estimator used in the Wald construction is also successfully employed, achieving the same statistical error as the Lasso based on the full sample, so long as $k = O(\sqrt{n/s^2 \log d})$.

Our contribution advances the understanding of distributed inference in the presence of large scale and distributed data, but there is still a great deal of work to be done in the area. We focus here on the fundamentals of hypothesis testing and estimation in the divide and conquer setting. Beyond this, there is a whole tool kit of statistical methodology designed for the single sample setting, whose split sample asymptotic properties are yet to be understood.

7. Proofs

In this section, we present the proofs of the main theorems appearing in Sections 3 and 4. The statements and proofs of several auxiliary lemmas appear in the Supplementary Material. To simplify notation, we take $\beta_v^H = 0$ without loss of generality.

7.1. Proofs for Section 3.1

The proof of Theorem 3.3, relies on the following lemma, which bounds the probability that optimization problems in (3.4) are feasible.

Lemma 7.1: Assume $\Sigma = \mathbb{E}(X_i X_i^T)$ satisfies $C_{\min} < \lambda_{\min}(\Sigma) < \lambda_{\max}(\Sigma) < C_{\max}$ as well as $\|\Sigma^{-1/2} X_1\|_{\psi_2} = \kappa$, then we have

$$\mathbb{P} \left(\max_{j=1, \dots, k} \|M^{(j)} \widehat{\Sigma}^{(j)} - I\|_{\max} \leq a \sqrt{\frac{\log d}{n}} \right) \geq 1 - 2kd^{-c_2},$$

where $c_2 = \frac{a^2 C_{\min}}{24e^2 \kappa^4 C_{\max}} - 2$.

Proof: The proof is an application of the union bound in Lemma 6.2 of Javanmard and Montanari (2014).

Proof of Theorem 3.3: For $1 \leq j \leq k$, let $\sqrt{n_k}(\widehat{\beta}^d(\mathcal{D}_j) - \beta^*) = Z^{(j)} + \Delta^{(j)}$, where $Z^{(j)} = \frac{1}{\sqrt{n_k}} M^{(j)} X^{(j)T} \epsilon^{(j)}$. From Theorem F.1, we know that as long as $\mathbf{m}_v^{(j)T} \widehat{\Sigma}^{(j)} \mathbf{m}_v^{(j)} \geq c > 0$ holds uniformly for $j = 1, \dots, d$,

$$\overline{\Delta} := \sqrt{n_k} \frac{1}{k} \sum_{j=1}^k \frac{\Delta_v^{(j)}}{\widehat{Q}^{(j)}} = o_{\mathbb{P}}(1).$$

Then we define

$$\bar{V}_n := \sqrt{n_k} \frac{1}{k} \sum_{j=1}^k \frac{Z_v^{(j)}}{\widehat{Q}^{(j)}} = \sum_{j=1}^k \frac{1}{k} \sum_{i \in \mathcal{I}_j} \xi_{iv}^{(j)}, \quad \text{where } \xi_{iv}^{(j)} := \frac{\mathbf{m}_v^{(j)T} X_i^{(j)} \epsilon_i^{(j)}}{\left(\mathbf{m}_v^{(j)T} \widehat{\Sigma}^{(j)} \mathbf{m}_v^{(j)} \right)^{1/2}}.$$

We now establish the asymptotic normality of \bar{V}_n by verifying the requirements of the Lindeberg-Feller central limit theorem (e.g. Kallenberg, 1997, Theorem 4.12). By the fact that ϵ_i is independent of X for all i and $\mathbb{E}[\epsilon_i] = 0$,

$$\begin{aligned} \mathbb{E}(\xi_{iv}^{(j)}) &= \mathbb{E}(\mathbb{E}(\xi_{iv}^{(j)} | X)) = \mathbb{E}\left\{\mathbb{E}\left[\mathbf{m}_v^{(j)T} \mathbf{X}_i^{(j)} \varepsilon_i^{(j)} / \left(n \mathbf{m}_v^{(j)T} \widehat{\Sigma}^{(j)} \mathbf{m}_v^{(j)}\right)^{1/2} \mid X\right]\right\} \\ &= \mathbb{E}\left\{\left(n \mathbf{m}_v^{(j)T} \widehat{\Sigma}^{(j)} \mathbf{m}_v^{(j)}\right)^{-1/2} \mathbf{m}_v^{(j)T} \mathbf{X}_i^{(j)} \mathbb{E}(\varepsilon_i^{(j)})\right\} = 0. \end{aligned}$$

By independence of $\{\varepsilon_i\}_{i=1}^n$ and the definition of $\widehat{\Sigma}^{(j)}$, we also have

$$\begin{aligned} \text{Var}(\bar{V}_n | X) &= \sum_{j=1}^k \sum_{i \in \mathcal{I}_j} \text{Var}(\xi_{iv}^{(j)} | X) \\ &= \frac{1}{n} \sum_{j=1}^k \left(\mathbf{m}_v^{(j)T} \widehat{\Sigma}^{(j)} \mathbf{m}_v^{(j)}\right)^{-1} \sum_{i \in \mathcal{I}_j} \mathbf{m}_v^{(j)T} \mathbf{X}_i^{(j)} \mathbf{X}_i^{(j)T} \mathbf{m}_v^{(j)} \text{Var}(\varepsilon_i^{(j)} | X) = \sigma^2. \end{aligned}$$

Therefore we have

$$\text{Var}(\bar{V}_n) = \mathbb{E}(\text{Var}(\bar{V}_n | X)) + \text{Var}(\mathbb{E}(\bar{V}_n | X)) = \sigma^2.$$

It only remains to verify the Lindeberg condition, i.e.,

$$\lim_{k \rightarrow \infty} \lim_{n_k \rightarrow \infty} \frac{1}{\sigma^2} \sum_{j=1}^k \sum_{i \in \mathcal{I}_j} \mathbb{E}\left[\left(\xi_{iv}^{(j)}\right)^2 \mathbb{1}\left\{\left|\xi_{iv}^{(j)}\right| > \varepsilon \sigma\right\}\right] = 0, \quad \forall \varepsilon > 0, \quad (7.1)$$

whose verification is relegated to the Appendix E of the Supplementary Material. Finally we reach the conclusion by the Slutsky's Theorem.

Proof of Corollary 3.5: Let $\mathcal{F}_n := \{\mathbf{m}_v^{(j)T} \widehat{\Sigma}^{(j)} \mathbf{m}_v^{(j)} \geq c > 0, j = 1, \dots, k\}$, where n is the total sample size. According to Theorem 3.3, when \mathcal{F}_n holds, we have

$$\lim_{n \rightarrow \infty} \mathbb{P}(\bar{\mathcal{S}}_n \leq t | \mathbf{X}) - \Phi(t) = 0.$$

From the proof of Lemma 13 in Javanmard and Montanari (2014), $\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{F}_n) = 1$. For any $t \in \mathbb{R}$ and $\delta > 0$, by applying dominating convergence Theorem to $\mathbb{1}\{|\mathbb{P}(\bar{\mathcal{S}}_n \leq t | \mathbf{X}) - \Phi(t)| > \delta \text{ and } \mathcal{F}_n \text{ holds}\}$, we have

$$\lim_{n \rightarrow \infty} \mathbb{P}(|\mathbb{P}(\bar{\mathcal{S}}_n \leq t | \mathbf{X}) - \Phi(t)| > \delta) = 0.$$

According to the dominate convergence theorem, since $\mathbb{P}(\bar{\mathcal{S}}_n \leq t | \mathbf{X}) \in [0, 1]$, we have

$$\lim_{n \rightarrow \infty} \frac{1}{n_k} \mathbb{P}(\bar{\mathcal{S}}_n \leq t) = \lim_{n \rightarrow \infty} \mathbb{E}[\mathbb{P}(\bar{\mathcal{S}}_n \leq t \mid \mathbf{X})] = \mathbb{E}[\lim_{n \rightarrow \infty} \mathbb{P}(\bar{\mathcal{S}}_n \leq t \mid \mathbf{X})] = \Phi(t).$$

Therefore, we complete the proof of the corollary.

The proofs of Theorem 3.8 and Corollary 3.9 are stated as an application of Lemmas E.7 and E.8 in the Supplementary Material, which apply under a more general set of requirements. We present the proof of Theorem 3.8 below and defer Corollary 3.9 to Appendix E in the Supplementary Materials.

Proof of Theorem 3.8: We verify (A1)–(A4) of Lemma E.7. For (A1), decompose the object of interest as

$$\frac{1}{n_k} \|X^{(j)} \widehat{\Theta}^{(j)}\|_{\max} = \frac{1}{n_k} \|X^{(j)} (\widehat{\Theta}^{(j)} - \Theta^*)\|_{\max} + \frac{1}{n_k} \|X^{(j)} \Theta^*\|_{\max} = \Delta_1 + \Delta_2,$$

where Δ_1 can be further decomposed and bounded by

$$\begin{aligned} \frac{1}{n_k} \|X^{(j)} (\widehat{\Theta}^{(j)} - \Theta^*)\|_{\max} &= \frac{1}{n_k} \max_{1 \leq i \leq n} \max_{1 \leq v \leq d} \left[|X_i^{(j)T} (\widehat{\Theta}_v^{(j)} - \Theta_v^*)| \right] \\ &\leq \frac{1}{n_k} \max_{1 \leq i \leq n} \|X_i\|_{\infty} \max_{1 \leq v \leq d} \|\widehat{\Theta}_v^{(j)} - \Theta_v^*\|_1. \end{aligned}$$

We have

$$\mathbb{P}(\Delta_1 > q/2) \leq \mathbb{P}\left(\max_{1 \leq v \leq d} \|\widehat{\Theta}_v^{(j)} - \Theta_v^*\|_1 > \frac{n_k q}{2} \right) < \psi$$

and by Condition 3.7, $\psi = \alpha(d^{-1}) = \alpha(k^{-1})$ for any $q \geq 2CMs_1(k/n)^{3/2} \cdot \sqrt{\log d}$, a fortiori for q a constant. Since X_j is sub-Gaussian, a matching probability bound can easily be obtained for Δ_2 , thus we obtain

$$\mathbb{P}(n_k^{-1} \|X^{(j)} \widehat{\Theta}^{(j)}\|_{\max} \leq 2\psi)$$

for $\psi = \alpha(k^{-1})$. (A2) and (A3) of Lemma E.7 are applications of Lemmas E.3 and E.4 respectively. To establish (A4), observe that

$$\left(\widehat{\Theta}_v^{(j)T} \nabla^2 \ell_{n_k}^{(j)}(\widehat{\beta}^{\lambda}(\mathcal{D}_j)) - e_v \right) = \Delta_1 + \Delta_2 + \Delta_3,$$

where $\Delta_1 = (\widehat{\Theta}_v^{(j)} - \Theta_v^*)^T \nabla^2 \ell_{n_k}^{(j)}(\widehat{\beta}^\lambda(\mathcal{D}_j))$, $\Delta_2 = \Theta_v^{*T} \left(\nabla^2 \ell_{n_k}^{(j)}(\widehat{\beta}^\lambda(\mathcal{D}_j)) - \nabla^2 \ell_{n_k}^{(j)}(\beta^*) \right)$ and $\Delta_3 = \Theta_v^{*T} \nabla^2 \ell_{n_k}^{(j)}(\beta^*) - e_v$. We thus consider $|\Delta_\ell(\widehat{\beta}^\lambda(\mathcal{D}_j) - \beta^*)|$ for $\ell = 1, 2, 3$.

$$\begin{aligned} |\Delta_2(\widehat{\beta}^\lambda(\mathcal{D}_j) - \beta^*)| &= \left| \frac{1}{n_k} \sum_{i \in \mathcal{J}_j} \Theta_v^{*T} X_i X_i^T (\widehat{\beta}^\lambda(\mathcal{D}_j) - \beta^*) [b''(X_i^T \widehat{\beta}^\lambda(\mathcal{D}_j)) - b''(X_i^T \beta^*)] \right| \\ &\leq U_3 \max_{1 \leq i \leq n} |\Theta_v^{*T} X_i| \frac{1}{n_k} \|X(\widehat{\beta}^\lambda(\mathcal{D}_j) - \beta^*)\|_2^2. \end{aligned}$$

$\mathbb{P} \left(\|X(\widehat{\beta}^\lambda(\mathcal{D}_j) - \beta^*)\|_2^2 \gtrsim n^{-1} s k \log(d/\delta) \right) < \delta$ by Lemma E.4, thus $\mathbb{P}(|\Delta_2(\widehat{\beta}^\lambda(\mathcal{D}_j) - \beta^*)| > t) < \delta$ for $t \asymp MU_3 n^{-1} s k \log(d/\delta)$. Invoking Hölder's inequality, Hoeffding's inequality and Condition 2.1, we also obtain, for $t \asymp n^{-1} s k \log(d/\delta)$,

$$\begin{aligned} &\mathbb{P} \left(|\Delta_3(\widehat{\beta}^\lambda(\mathcal{D}_j) - \beta^*)| > t \right) \\ &\leq \mathbb{P} \left(\left\| \Theta_v^{*T} \left(\frac{1}{n_k} \sum_{i \in \mathcal{J}_j} b''(X_i^T \beta^*) X_i X_i^T \right) - e_v \right\|_{\max} \|\widehat{\beta}^\lambda(\mathcal{D}_j) - \beta^*\|_1 > t \right). \end{aligned}$$

Therefore $\mathbb{P}(|\Delta_2(\widehat{\beta}^\lambda(\mathcal{D}_j) - \beta^*)| > t) < 2\delta$. Finally, with $t = n^{-1}(s \vee s_1)k \log(d/\delta)$,

$$\mathbb{P} \left(|\Delta_1(\widehat{\beta}^\lambda(\mathcal{D}_j) - \beta^*)| > t \right) \leq \mathbb{P} \left(\frac{1}{n_k} \left\| \sum_{i \in \mathcal{J}_j} X_i^T (\widehat{\Theta}_v - \Theta_v) b''(X_i^T \widehat{\beta}^\lambda(\mathcal{D}_j)) \right\|_2 \left\| X^{(j)T} (\widehat{\beta}^\lambda(\mathcal{D}_j) - \beta^*) \right\|_2 > t \right),$$

hence $\mathbb{P}(|\Delta_1(\widehat{\beta}^\lambda(\mathcal{D}_j) - \beta^*)| > t) < 2\delta$. This follows because by Lemma E.4,

$$\mathbb{P} \left(\left\| \frac{1}{n_k} X^{(j)} (\widehat{\beta}^\lambda(\mathcal{D}_j) - \beta^*) \right\|_2 \gtrsim n^{-1/2} \sqrt{s k \log(d/\delta)} < \delta \right)$$

and by Lemma C.4 of Ning and Liu (2014),

$$\mathbb{P} \left(\left\| \frac{1}{n_k} \sum_{i \in \mathcal{J}_j} X_i^T (\widehat{\Theta}_v - \Theta_v) b''(X_i^T \widehat{\beta}^\lambda(\mathcal{D}_j)) \right\|_2 \gtrsim n^{-1/2} \sqrt{s_1 k \log(d/\delta)} < \delta \right)$$

7.2. Proofs for Theorems in Section 3.2

The proof of Theorem 3.11 relies on several preliminary lemmas, collected in Appendix E in the Supplementary Material. Without loss of generality we set $H_0: \beta_v^* = 0$ to ease notation.

Proof of Theorem 3.11: Since $\bar{S}(0) = k^{-1} \sum_{j=1}^k \hat{S}^{(j)}(0, \hat{\beta}_{-v}^\lambda(\mathcal{D}_j))$, and (B1)–(B4) of Condition E.9 in the Supplementary Material are fulfilled under Conditions 3.6 and 2.1 by Lemma E.10 (see Appendix E in the Supplementary Material). The proof is now simply an application of Lemma E.13 in the Supplementary Material with $\beta_v^* = 0$ under the restriction of the null hypothesis.

Proof of Lemma 3.14: The proof is an application of Lemma E.16 in the Supplementary Material, noting that (B1)–(B5) of Condition E.9 in the Supplementary Material are fulfilled under Conditions 3.6 and 2.1 by Lemmas E.10 and E.11 in the Supplementary Material.

7.3. Proofs for Theorems in Section 4

Recall from Section 2 that for an arbitrary matrix M , M_{ℓ} denotes the transposed ℓ^{th} row of M and $[M]_{\ell}$ denotes the ℓ^{th} column of M .

Proof of Theorem 4.3: By Lemma 4.1 and $k = O(\sqrt{nl/(s^2 \log d)})$, there exists a sufficiently large C_0 such that for the event $\mathcal{E} = \{\|\bar{\beta}^d - \beta^*\|_{\infty} \leq C_0 \sqrt{\log d/n}\}$, we have $\mathbb{P}(\mathcal{E}) = 1 - c/d$. We choose $\nu = C_0 \sqrt{\log d/n}$, which implies that, under \mathcal{E} , we have $\nu = \|\bar{\beta}^d - \beta^*\|_{\infty}$.

Let \mathcal{S} be the support of β^* . The derivations in the remainder of the proof hold on the event \mathcal{E} . Observe $\mathcal{T}_{\nu}(\bar{\beta}_{\mathcal{S}^c}^d) = \mathbf{0}$ as $\|\bar{\beta}_{\mathcal{S}^c}^d\|_{\infty} \leq \nu$. For $j \in \mathcal{S}$, if $|\beta_j^*| \geq 2\nu$, we have

$|\bar{\beta}_j^d| \geq |\beta_j^*| - \nu \geq \nu$ and thus $|\mathcal{T}_{\nu}(\bar{\beta}_j^d) - \beta_j^*| = |\bar{\beta}_j^d - \beta_j^*| \leq \nu$. While if $|\beta_j^*| < 2\nu$, $|\mathcal{T}_{\nu}(\bar{\beta}_j^d) - \beta_j^*| \leq |\beta_j^*| \vee |\bar{\beta}_j^d - \beta_j^*| \leq 2\nu$. Therefore, on the event \mathcal{E} ,

$$\begin{aligned} \|\mathcal{T}_{\nu}(\bar{\beta}^d) - \beta^*\|_2 &= \|\mathcal{T}_{\nu}(\bar{\beta}_{\mathcal{S}}^d) - \beta_{\mathcal{S}}^*\|_2 \leq 2\sqrt{s\nu} \\ \text{and } \|\mathcal{T}_{\nu}(\bar{\beta}^d) - \beta^*\|_{\infty} &= \|\mathcal{T}_{\nu}(\bar{\beta}_{\mathcal{S}}^d) - \beta_{\mathcal{S}}^*\|_{\infty} \leq 2\nu. \end{aligned}$$

The statement of the theorem follows because $\nu = C_0 \sqrt{\log d/n}$ and $\mathbb{P}(\mathcal{E}) = 1 - c/d$. Following the same reasoning, on the event

$\mathcal{E}': = \mathcal{E} \cup \{\|\bar{\beta}^d - \beta^*\|_{\infty} \leq C_0 \sqrt{\log d/n}\} \cup \{\|\hat{\beta}^d - \bar{\beta}^d\|_{\infty} \leq C_0 s k \log d/n\}$, we have

$$\begin{aligned} \left\| \mathcal{T}_\nu(\tilde{\beta}^d) - \mathcal{T}_\nu(\hat{\beta}_S^d) \right\|_2 &= \left\| \mathcal{T}_\nu(\tilde{\beta}_S^d) - \mathcal{T}_\nu(\hat{\beta}_S^d) \right\|_2 \\ &\leq \left\| \tilde{\beta}_S^d - \hat{\beta}_S^d \right\|_2 \leq \sqrt{s} \left\| \tilde{\beta}_S^d - \hat{\beta}_S^d \right\|_\infty \leq C s^{3/2} k \log d/n. \end{aligned}$$

As Lemma 4.1 also gives $\mathbb{P}(\mathcal{E}') = 1 - c/d$, the proof is complete.

Proof of Corollary 4.10: By an analogous proof strategy to that of Theorem 4.8,

$$\begin{aligned} \left| [\mathcal{T}_\zeta(\bar{\Theta})]_{\nu\nu} - \Theta_{\nu\nu}^* \right| &= O_p\left(\sqrt{n^{-1} \log d}\right) = o_{\mathbb{P}}(1) \text{ under the conditions of the Corollary provided} \\ k &= o\left((s \vee s_1) \log d\right)^{-1} \sqrt{n}. \end{aligned}$$

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

The authors thank Weichen Wang, Jason Lee and Yuekai Sun for helpful comments.

References

- Bickel PJ. One-step huber estimates in the linear model. *Journal of the American Statistical Association*. 1975; 70:428–434.
- Bühlmann P, van de Geer S. *Statistics for high-dimensional data: methods, theory and applications* Springer; 2011
- Candes E, Tao T. The Dantzig selector: statistical estimation when p is much larger than n . *Ann Statist*. 2007; 35:2313–2351.
- Chen X, Xie M. Tech Rep 2012-01 Department of Statistics, Rutgers University; 2012 A split and conquer approach for analysis of extraordinarily large data.
- Chernozhukov V, Chetverikov D, Kato K. Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors. *Ann Statist*. 2013; 41:2786–2819.
- Cox DR, Hinkley DV. *Theoretical statistics* Chapman and Hall; London: 1974
- Fan J, Guo S, Hao N. Variance estimation using refitted cross-validation in ultrahigh dimensional regression. *J R Stat Soc Ser B Stat Methodol*. 2012; 74:37–65.
- Fan J, Han F, Liu H. Challenges of big data analysis. *National Sci Rev*. 2014; 1:293–314.
- Fan J, Li R. Variable selection via nonconcave penalized likelihood and its oracle properties. *J Amer Statist Assoc*. 2001; 96:1348–1360.
- Fan J, Lv J. Nonconcave penalized likelihood with np-dimensionality. *Information Theory, IEEE Transactions on*. 2011; 57:5467–5484.
- Fan J, Song R. Sure independence screening in generalized linear models with NP-dimensionality. *Ann Statist*. 2010; 38:3567–3604.
- Javanmard A, Montanari A. Confidence intervals and hypothesis testing for high-dimensional regression. *Journal of Machine Learning Research*. 2014; 15:2869–2909.
- Kallenberg O. *Probability and its Applications* (New York) Springer-Verlag; New York: 1997 Foundations of modern probability.
- Kleiner A, Talwalkar A, Sarkar P, Jordan MI. A scalable bootstrap for massive data. *J R Stat Soc Ser B Stat Methodol*. 2014; 76:795–816.
- Lee JD, Sun Y, Liu Q, Taylor JE. Communication-efficient sparse regression: a one-shot approach. 2015 ArXiv 1503.04337.

- Liu Q, Ihler AT. Distributed estimation, information loss and exponential families. *Advances in Neural Information Processing Systems*. 2014
- Loh P-L, Wainwright MJ. Regularized m -estimators with nonconvexity: Statistical and algorithmic theory for local optima. In: Burges C, Bottou L, Welling M, Ghahramani Z, Weinberger K, editors *Advances in Neural Information Processing Systems 2013* 26476484
- Loh P-L, Wainwright MJ. Regularized M -estimators with nonconvexity: statistical and algorithmic theory for local optima. *J Mach Learn Res*. 2015; 16:559–616.
- Meinshausen N, Bühlmann P. High-dimensional graphs and variable selection with the lasso. *Ann Statist*. 2006; 34:1436–1462.
- Negahban S, Yu B, Wainwright MJ, Ravikumar PK. A unified framework for high-dimensional analysis of m -estimators with decomposable regularizers. *Advances in Neural Information Processing Systems*. 2009
- Ning Y, Liu H. A General Theory of Hypothesis Tests and Confidence Regions for Sparse High Dimensional Models. 2014 ArXiv 1412.8765.
- Rosenblatt JD, Nadler B. On the optimality of averaging in distributed statistical learning. *Information and Inference*. 2016:iaw013.
- Tibshirani R. Regression shrinkage and selection via the lasso. *J Roy Statist Soc Ser B*. 1996; 58:267–288.
- van de Geer S, Bühlmann P, Ritov Y, Dezeure R. On asymptotically optimal confidence regions and tests for high-dimensional models. *Ann Statist*. 2014; 42:1166–1202.
- Vershynin R. Introduction to the non-asymptotic analysis of random matrices. 2010 arXiv preprint arXiv:1011.3027.
- Wang Z, Liu H, Zhang T. Optimal computational and statistical rates of convergence for sparse nonconvex learning problems. *Ann Statist*. 2014a; 42:2164–2201.
- Wang Z, Liu H, Zhang T. Optimal computational and statistical rates of convergence for sparse nonconvex learning problems. *Ann Statist*. 2014b; 42:2164–2201.
- Zhang C-H. Nearly unbiased variable selection under minimax concave penalty. *Ann Statist*. 2010; 38:894–942.
- Zhang C-H, Zhang SS. Confidence intervals for low dimensional parameters in high dimensional linear models. *J R Stat Soc Ser B Stat Methodol*. 2014; 76:217–242.
- Zhang C-H, Zhang T. A general theory of concave regularization for high-dimensional sparse estimation problems. *Statistical Science*. 2012; 27:576–593.
- Zhang Y, Duchi JC, Wainwright MJ. Divide and Conquer Kernel Ridge Regression: A Distributed Algorithm with Minimax Optimal Rates. 2013 ArXiv e-prints.
- Zhao T, Cheng G, Liu H. A Partially Linear Framework for Massive Heterogeneous Data. 2014a ArXiv 1410.8570.
- Zhao T, Kolar M, Liu H. A General Framework for Robust Testing and Confidence Regions in High-Dimensional Quantile Regression. 2014b ArXiv 1412.8724.

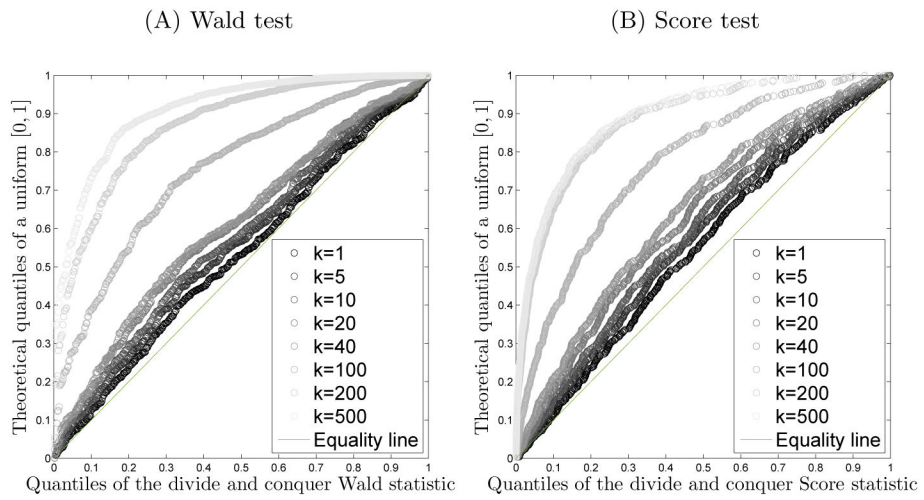


Fig 1. QQ plots of the p -values of the Wald (A) and score (B) divide and conquer test statistics against the theoretical quantiles of the uniform $[0,1]$ distribution under the null hypothesis.

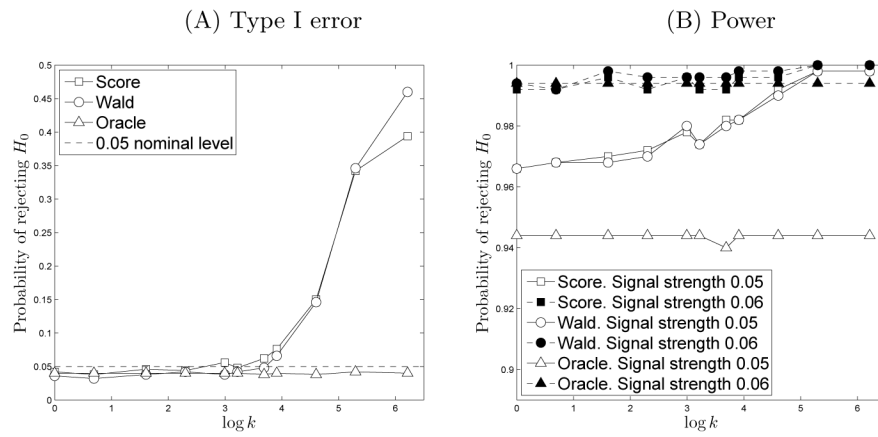


Fig 2. (A) Estimated probabilities of type I error for the Wald and score tests as a function of k . (B) Estimated power with signal strength 0.05 and 0.06 for the Wald, and score tests as a function of k .

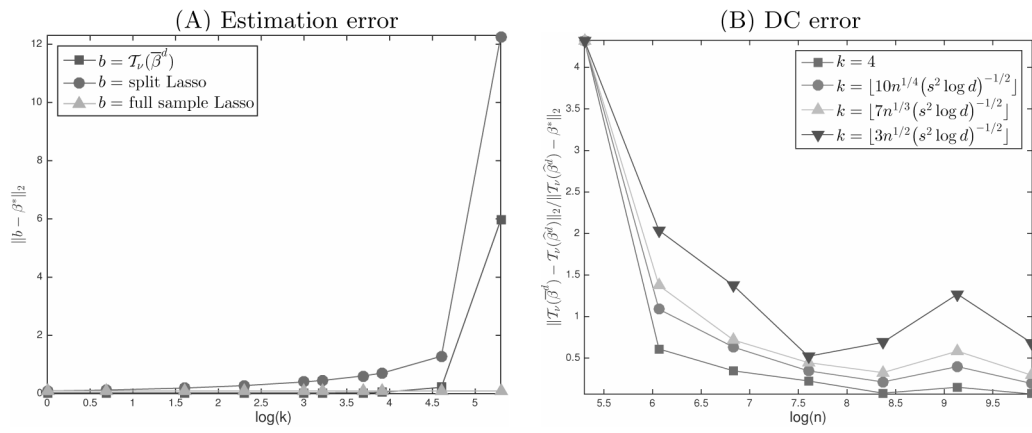


Fig 3. (A): Statistical error of the DC estimator, split Lasso and the full sample Lasso for $k \in \{1, 2, 5, 10, 20, 25, 40, 50, 100, 200\}$ when $n = 5000, d = 5000$. (B): Euclidean norm difference between the DC thresholded debiased estimator and its full sample analogue.

Computation time for the divide and conquer testing and estimation, where $k = 1$ corresponds to the non-splitting case and $k > 1$ corresponds to the distributed case.

Table 1

k	1	2	5	10	20	25	50	100	200
Score test (\$)	364.39	73.22	35.09	23.61	23.56	20.78	24.13	37.53	64.67
Wald test (\$)	426.23	68.95	19.66	10.09	6.70	5.71	3.88	2.60	1.91
$\mathcal{J}_{\sqrt{\lambda}\beta^k}(10^3s)$	61.50	30.00	7.92	6.58	4.48	2.94	2.64	2.11	1.66
Split Lasso (\$)	89.18	32.02	34.57	6.47	4.87	4.16	2.56	1.92	2.64