# Communication-Efficient Accurate Statistical Estimation\*

Jianqing Fan, Yongyi Guo and Kaizheng Wang

#### **Abstract**

When the data are stored in a distributed manner, direct application of traditional statistical inference procedures is often prohibitive due to communication cost and privacy concerns. This paper develops and investigates two Communication-Efficient Accurate Statistical Estimators (CEASE), implemented through iterative algorithms for distributed optimization. In each iteration, node machines carry out computation in parallel and communicate with the central processor, which then broadcasts aggregated information to node machines for new updates. The algorithms adapt to the similarity among loss functions on node machines, and converge rapidly when each node machine has large enough sample size. Moreover, they do not require good initialization and enjoy linear converge guarantees under general conditions. The contraction rate of optimization errors is presented explicitly, with dependence on the local sample size unveiled. In addition, the improved statistical accuracy per iteration is derived. By regarding the proposed method as a multi-step statistical estimator, we show that statistical efficiency can be achieved in finite steps in typical statistical applications. In addition, we give the conditions under which the one-step CEASE estimator is statistically efficient. Extensive numerical experiments on both synthetic and real data validate the theoretical results and demonstrate the superior performance of our algorithms.

<sup>\*</sup>Supported by NSF grants DMS-1662139 and DMS-1712591, NIH grant 2R01-GM072611-14, and ONR grant N00014-19-1-2120.

## 1 Introduction

Statistical inference in modern era faces tremendous challenge on computation and storage. The exceedingly large size of data often makes it impossible to store all of them on a single machine. Moreover, many applications have individual agents (e.g. local governments, research labs, hospitals, smart phones) collecting data independently. Communication between them is prohibitively expensive due to the limited bandwidth, and direct data sharing has also raised privacy and lost of ownership concerns. These constraints make it necessary to develop methodologies for distributed systems, solving statistical problems with divide-and-conquer procedures and communicating only certain summary statistics. In modern distributed computing architectures, the speeds of intra-node computation and inter-node communication may differ by a factor of 1000 (Lan et al., 2018). It is then desirable to conduct expensive computation on node machines and communicate as few rounds as possible.

Distributed statistical inference has received considerable attention in the past few years, covering a wide spectrum of topics including M-estimation (Zhang et al., 2013; Chen and Xie, 2014; Shamir et al., 2014; Wang et al., 2017a; Lee et al., 2017b; Battey et al., 2018; Jordan et al., 2018; Wang et al., 2018a; Shi et al., 2018; Chen et al., 2018; Banerjee et al., 2019), principal component analysis (Fan et al., 2017; Garber et al., 2017), nonparametric regression (Zhang et al., 2015; Shang and Cheng, 2017; Szabo and van Zanten, 2017; Han et al., 2018), quantile regression (Volgushev et al., 2017; Chen et al., 2018), bootstrap (Kleiner et al., 2014), confidence intervals (Jordan et al., 2018; Wang et al., 2018b; Chen et al., 2018), Bayesian methods (Suchard et al., 2010; Wang and Dunson, 2013; Jordan et al., 2018), and so on. In the commonly-used setting, the overall dataset is partitioned and stored on m node machines, which are connected to a central processor. Most of the approaches studied in this literature only require one round of communication: the node machines conduct inference in parallel and send their results to the central processor, which then aggregates the information and outputs a final result. As typical examples, Zhang et al. (2013) average the M-estimators obtained by node machines; Battey et al. (2018) average debiased estimators; and Fan et al. (2017) define an average for subspaces and compute it via eigen-decomposition. While these one-shot methods are communication-efficient, they only work with a small number of node machines (e.g.  $m = o(\sqrt{N})$ , where N is the total sample size) and require large sample size on each of them, as their theories heavily rely on asymptotic expansions of certain estimators (Zhang et al., 2013; Rosenblatt and Nadler, 2016). Since the conditions are easily violated, their performance may well be sub-optimal.

Multi-round procedures come as a remedy for this, where local computation and global aggregation are repeatedly performed. On the one hand, the central processor gathers and broadcasts overall information for node machines to improve their estimates accordingly. On the other hand, the similar structures of data on node machines as well as their computational power are exploited by the algorithm. It is then possible to achieve optimal statistical precision after a few rounds of communication, under broader settings than those for oneshot procedures. Shamir et al. (2014) proposes a Distributed Approximate NEwton (DANE) algorithm where, in each iteration, each node machine minimizes a modified loss function based on its own samples and the gradient information from all other machines obtained through communication. However, for non-quadratic losses, the analysis in Shamir et al. (2014) does not imply any advantage of DANE in terms of communication over distributed implementation of gradient descent. Other approximate Newton algorithms include Zhang and Xiao (2015), Mahajan et al. (2015), Wang et al. (2018a) and Crane and Roosta (2019). A recent work by Chen et al. (2018) approximate the Newton step by first-order stochastic algorithms. In addition, Jordan et al. (2018) develops a Communication-efficient Surrogate Likelihood (CSL) framework for estimation and inference in regular parametric models, highdimensional penalized regression, and Bayesian statistics. A similar method for penalized regression also appear independently in Wang et al. (2017a). These methods no longer have restrictions on the number of machines such as  $m = o(\sqrt{N})$ .

Due to the nature of Newton-type methods, existing theories for these algorithms heavily rely on good initialization or even self-concordance assumption on loss functions. They essentially focus on improving an initial estimator that is already consistent but not efficient, whose ideas coincide with the classical one-step estimator (Bickel, 1975). Such initialization itself needs additional efforts and assumptions. Moreover, current results still require each machine to have sufficiently many samples so that loss functions on different machines are similar to each other. These all make the proposed methods unreliable in practice.

Aside from distributed statistical inference, there has also been a vast literature in distributed optimization since the pioneering works in the 80s (Bertsekas and Tsitsiklis, 1989). The ADMM (Boyd et al., 2011) is a celebrated example among the numerous algorithms designed to handle deterministic optimization problems with minimum structural assumption. While having theoretical guarantees under general conditions (Deng and Yin, 2016; Hong and Luo, 2017), the convergence can be quite slow. Recent developments in distributed optimization include dual algorithms (Yang, 2013; Jaggi et al., 2014; Smith et al., 2018), algorithms for feature-partitioned problems Recht et al. (2011); Richtárik and Takáč (2016),

resampling-based algorithms (Lee et al., 2017a; Wang et al., 2017b), decentralized algorithms (Duchi et al., 2012; Lan et al., 2018), federated optimization (Konečný et al., 2015; Chen et al., 2017), communication complexity (Arjevani and Shamir, 2015; Woodworth et al., 2018), to name a few. This list is by no means exhaustive for this booming area. In the statistical setting we are interested in, most of the algorithms above cannot fully utilize the similarity among loss functions on node machines.

In this paper, we develop and study two Communication-Efficient Accurate Statistical Estimators (CEASE) based on multi-round algorithms for distributed statistical estimation. The samples are stored on m node machines connected to a central processor. For simplicity, we assume that all the node machines have the same sample size n. Each node machine has a regularized empirical risk function  $f_k + g$  defined by the samples stored there, and the goal is to compute the minimizer of the overall regularized risk function  $\frac{1}{m}\sum_{k=1}^{m}f_k+g$  to statistical precision. The algorithms alternate between computation on node machines and aggregation on the central processor in a communication-efficient way. When n is sufficiently large, their rates of convergence are better than or comparable to existing statistical methods designed for this large-sample regime. Even for moderate or small n, they are still guaranteed to converge linearly even in the absence of good initializations, while other statistical methods fail. In addition, our algorithms take advantage of the similarity among  $\{f_k\}_{k=1}^m$  in statistical applications, and thus improve over general-purpose distributed optimization algorithms like ADMM. To some extent, they interpolate between distributed algorithms for statistical estimation and general deterministic problems. Theoretical findings are verified by extensive numerical experiments.

From a technical point of view, our algorithms use the proximal point algorithm (Rockafellar, 1976; Parikh and Boyd, 2014) as the backbone and obtain inexact one-step updates in a distributed manner. This turns out to be crucial for proving convergence under general conditions, without good initialization or large sample size n on each node machine. Moreover, it makes our algorithms reliable in practice. Our perspective and techniques are potentially useful for analyzing other distributed optimization algorithms.

The rest of this paper is organized as follows. Section 2 introduces the problem setup and presents two vanilla algorithms for the large-sample regime. Section 3 proposes two advanced algorithms and analyzes their theoretical properties under general conditions. Section 4 uses numerical experiments on both synthetic and real data to validate the theoretical results. Section 5 finally concludes the paper and discusses possible future directions.

#### **Notations**

Here we list notations to be used throughout the paper. [n] denotes the set  $\{1, 2, \dots, n\}$  for any positive integer n. For two sequences  $\{a_n\}_{n=1}^{\infty}$  and  $\{b_n\}_{n=1}^{\infty}$ , we write  $a_n = O(b_n)$  or  $a_n \lesssim b_n$  if there exists a constant C > 0 such that  $a_n \leq Cb_n$  holds for sufficiently large n; and  $a_n \times b_n$  if  $a_n = O(b_n)$  and  $b_n = O(a_n)$ . Given a Euclidean space  $\mathbb{R}^k$  where k is clear from the context,  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^k$ , and r > 0, we define  $B(\mathbf{x}, r) = \{\mathbf{z} \in \mathbb{R}^k : \|\mathbf{z} - \mathbf{x}\|_2 \leq r\}$  to be a closed ball and  $\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{j=1}^k x_j y_j$  to be the inner product. For a convex function h on  $\mathbb{R}^k$ , we let  $\partial h(\mathbf{x})$  be its sub-differential set at  $\mathbf{x} \in \mathbb{R}^k$ , and  $\operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^k} h(\mathbf{x})$  be the set of its minimizers if  $\inf_{\mathbf{x} \in \mathbb{R}^k} h(\mathbf{x}) > -\infty$ . We use  $\|\cdot\|_2$  to denote the  $\ell_2$  norm of a vector or operator norm of a matrix. For two sequences of random variables  $\{X_n\}_{n=1}^{\infty}$  and  $\{Y_n\}_{n=1}^{\infty}$  where  $Y_n \geq 0$ , we write  $X_n = O_{\mathbb{P}}(Y_n)$  if for any  $\varepsilon > 0$  there exists C > 0 such that  $\mathbb{P}(|X_n| \geq CY_n) \leq \varepsilon$  for sufficiently large n. We use  $\|X\|_{\psi_2} = \sup_{p\geq 1} \mathbb{E}^{1/p} |X|^p$  to refer to the sub-Gaussian norm of random variable X, and  $\|\mathbf{X}\|_{\psi_2} = \sup_{\|u\|_2=1} \|\langle u, \mathbf{X} \rangle\|_{\psi_2}$  to denote the sub-Gaussian norm of random vector  $\mathbf{X}$ .

# 2 Distributed estimation in large sample regimes

## 2.1 Problem setup

Suppose there is an unknown probability distribution  $\mathcal{P}$  over some sample space  $\mathcal{X}$ . For any parameter  $\boldsymbol{\theta} \in \mathbb{R}^p$ , define its population risk  $F(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{X} \sim \mathcal{P}} \ell(\boldsymbol{\theta}; \mathbf{X})$  based on a loss function  $\ell : \mathbb{R}^p \times \mathcal{X} \to \mathbb{R}$ . In parametric inference problems,  $\ell$  is often chosen as the negative log-likelihood function of some parametric family. Under mild conditions, F is well-defined and has a unique minimizer  $\boldsymbol{\theta}^*$ . A ubiquitous problem in statistics and machine learning is to estimate  $\boldsymbol{\theta}^*$  given i.i.d. samples  $\{\mathbf{X}_i\}_{i=1}^N$  from  $\mathcal{P}$ , and the minimizer of the empirical risk  $f(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \ell(\boldsymbol{\theta}; \mathbf{X}_i)$  becomes a natural candidate. To achieve desirable precision in high-dimensional problems, it is often necessary to incorporate prior knowledge of  $\boldsymbol{\theta}^*$  into the estimation procedure. To this end, the regularized empirical risk minimization

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^p} \left\{ f(\boldsymbol{\theta}) + g(\boldsymbol{\theta}) \right\}, \tag{2.1}$$

provides a principled approach, where  $g(\boldsymbol{\theta})$  is a deterministic panelty function. Common choices for  $g(\boldsymbol{\theta})$  include the  $\ell_2$  penalty  $\lambda \|\boldsymbol{\theta}\|_2^2$  (Hoerl and Kennard, 1970), the  $\ell_1$  penalty  $\lambda \|\boldsymbol{\theta}\|_1$  (Tibshirani, 1996), and a family of folded concave penalty functions  $\|p_{\lambda}(|\boldsymbol{\theta}|)\|_1$  such

as SCAD (Fan and Li, 2001) and MCP (Zhang et al., 2010), where  $\lambda > 0$  is a regularization parameter. Throughout the paper, we assume that both  $\ell$  and g are convex in  $\theta$ , and  $\ell$  is twice continuously differentiable in  $\theta$ . We allow g to be non-smooth (e.g. the  $\ell_1$  penalty).

Consider the distributed setting where the samples  $\{\mathbf{X}_i\}_{i=1}^N$  are stored on m machines connected to a central processor. Define  $\mathcal{I}_k$  to be the indices of samples on the k-th machine, and  $f_k(\boldsymbol{\theta}) = \frac{1}{|\mathcal{I}_k|} \sum_{i \in \mathcal{I}_k} \ell(\boldsymbol{\theta}; \mathbf{X}_i)$  the associated empirical loss. For simplicity, we assume that  $\{\mathcal{I}_k\}_{k=1}^m$  are disjoint, N is a multiple of m, and  $|\mathcal{I}_k| = n = N/m$  for all  $k \in [m]$ . Then (2.1) can be rewritten as

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^p} \left\{ f(\boldsymbol{\theta}) + g(\boldsymbol{\theta}) \right\}, \qquad f(\boldsymbol{\theta}) = \frac{1}{m} \sum_{k=1}^m f_k(\boldsymbol{\theta}). \tag{2.2}$$

Each machine k only has access to its local data and hence local loss function  $f_k$  and the penalty g. We aim to solve (2.2) in a distributed manner with both statistical efficiency and communication-efficiency.

## 2.2 Adaptive gradient enhancements and distributed algorithms

We drop the regularization term for now and consider the empirical risk minimization problem  $\min_{\boldsymbol{\theta} \in \mathbb{R}^p} f(\boldsymbol{\theta})$  for estimating  $\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^p} F(\boldsymbol{\theta})$ . In some problems, direct minimization of f is costly, while it is easy to obtain some rough estimate  $\bar{\boldsymbol{\theta}}$  that is close to  $\boldsymbol{\theta}^*$  but not as accurate as the global minimizer  $\hat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^p} f(\boldsymbol{\theta})$ . Bickel (1975) proposes the one-step estimator based on the local quadratic approximation and shows that it is as efficient as  $\hat{\boldsymbol{\theta}}$  if the initial estimator  $\bar{\boldsymbol{\theta}}$  is accuracy enough. Iterating this further results in multiple-step estimators that improve the optimization error and hence statistical errors when the initial estimator is not good enough (Robinson, 1988). This inspires us to refine an existing estimator using some proxy of f.

In the distributed environment, starting from an initial estimator  $\bar{\boldsymbol{\theta}}$ , only the gradient vector  $\nabla f(\bar{\boldsymbol{\theta}})$  can easily be communicated and hence the linear function  $f^{(1)}(\boldsymbol{\theta}) = f(\bar{\boldsymbol{\theta}}) + \langle \nabla f(\bar{\boldsymbol{\theta}}), \boldsymbol{\theta} - \bar{\boldsymbol{\theta}} \rangle$ , the first-order Taylor expansion of f around  $\bar{\boldsymbol{\theta}}$ . The object function to be minimized can be written as

$$f(\boldsymbol{\theta}) = f^{(1)}(\boldsymbol{\theta}) + R(\boldsymbol{\theta}), \quad \text{where} \quad R(\boldsymbol{\theta}) = f(\boldsymbol{\theta}) - f^{(1)}(\boldsymbol{\theta}).$$

Since the linear function  $f^{(1)}(\boldsymbol{\theta})$  can easily be communicated to each node machine whereas

 $R(\cdot)$  can not, the latter is naturally replaced by its subsampled version at node k:

$$R_k(\boldsymbol{\theta}) = f_k(\boldsymbol{\theta}) - [f_k(\bar{\boldsymbol{\theta}}) + \langle \nabla f_k(\bar{\boldsymbol{\theta}}), \boldsymbol{\theta} - \bar{\boldsymbol{\theta}} \rangle],$$

where  $f_k(\boldsymbol{\theta})$  is the loss function based on the data at node k. With this replacement, the target of optimization at node k becomes  $f^{(1)}(\boldsymbol{\theta}) + R_k(\boldsymbol{\theta})$ , which equals to

$$f_k(\boldsymbol{\theta}) - \langle \nabla f_k(\bar{\boldsymbol{\theta}}) - \nabla f(\bar{\boldsymbol{\theta}}), \boldsymbol{\theta} \rangle$$

up to an additive constant. This function will be called gradient-enhenced loss (GEL) function, in which the gradient at point  $\bar{\theta}$  based on the local data is replaced by the global one. This function has one very nice fixed point at the global minimum  $\hat{\theta}$ : the minimizer of the adaptive gradient-enhanced function at  $\bar{\theta} = \hat{\theta}$  is still  $\hat{\theta}$ . This can easily be seen by computing the gradient at the point  $\hat{\theta}$ .

The idea of using such an adaptive enhanced function has been proposed in Shamir et al. (2014) and Jordan et al. (2018), though the motivations are different. Jordan et al. (2018) develop a Communication-efficient Surrogate Likelihood (CSL) method using the GEL function  $f_1(\theta) - \langle \nabla f_1(\bar{\theta}) - \nabla f(\bar{\theta}), \theta \rangle$  on the first machine, uses the minimizer on that machine as a new estimate, and iterates these steps until convergence. In the presence of a regularizer g in (2.1), one simply adds g to the gradient-enhanced loss; see the Algorithm 1 below. It is also studied by Shamir et al. (2014) and Wang et al. (2017a) under certain settings.

#### Algorithm 1 CSL (Jordan et al., 2018)

**Input**: Initial value  $\theta_0$ , number of iterations T.

For  $t = 0, 1, 2, \dots, T - 1$ :

- Each machine evaluates  $\nabla f_k(\boldsymbol{\theta}_t)$  and sends to the 1<sup>st</sup> machine;
- The 1<sup>st</sup> machine computes  $\nabla f(\boldsymbol{\theta}_t) = \frac{1}{m} \sum_{k=1}^m \nabla f_k(\boldsymbol{\theta}_t)$  and

$$\boldsymbol{\theta}_{t+1} = \operatorname{argmin}_{\boldsymbol{\theta}} \left\{ f_1(\boldsymbol{\theta}) + g(\boldsymbol{\theta}) - \langle \nabla f_1(\boldsymbol{\theta}_t) - \nabla f(\boldsymbol{\theta}_t), \boldsymbol{\theta} \rangle \right\}$$

and broadcasts to other machines.

#### Output: $\theta_T$ .

Note that in Algorithm 1, only the first machine solves optimization problems and others just evaluate gradients. These machines are idling while the first one is working hard. To

fully utilize the computing power of machines and accelerate convergence, all the machines can optimize their corresponding GEL functions in parallel and the central processor then aggregates the results. This is motivated by the Distributed Approximate NEwton (DANE) algorithm (Shamir et al., 2014). Algorithm 2 describes the procedure in detail. Intuitively, the averaging step requires little computation but helps reduce the variance of estimators on node machines and enhance the accuracy.

#### Algorithm 2 GEL Method

**Input**: Initial value  $\theta_0$ , number of iterations T.

For  $t = 0, 1, 2, \dots, T - 1$ :

- Each machine evaluates  $\nabla f_k(\boldsymbol{\theta}_t)$  and sends to the central processor;
- The central processor computes  $\nabla f(\boldsymbol{\theta}_t) = \frac{1}{m} \sum_{k=1}^m \nabla f_k(\boldsymbol{\theta}_t)$  and broadcasts to machines;
- Each machine computes

$$\boldsymbol{\theta}_{t,k} = \operatorname{argmin}_{\boldsymbol{\theta}} \left\{ f_k(\boldsymbol{\theta}) + g(\boldsymbol{\theta}) - \langle \nabla f_k(\boldsymbol{\theta}_t) - \nabla f(\boldsymbol{\theta}_t), \boldsymbol{\theta} \rangle \right\}$$

and sends to the central processor;

• The central processor computes  $\boldsymbol{\theta}_{t+1} = \frac{1}{m} \sum_{k=1}^{m} \boldsymbol{\theta}_{t,k}$  and broadcasts to machines.

Output:  $\theta_T$ .

# 2.3 Contracting optimization errors

In this subsection, we first present deterministic (almost sure) results for Algorithms 1 and 2 based on high-level structural assumptions. We will then apply the results to the statistical setting in the next subsection.

**Definition 2.1.** Let  $h : \mathbb{R}^p \to \mathbb{R}$  be a convex function,  $\Omega \subseteq \mathbb{R}^p$  be a convex set, and  $\rho \geq 0$ . h is said to be  $\rho$ -strongly convex in  $\Omega$  if  $h(\mathbf{y}) \geq h(\mathbf{x}) + \langle \mathbf{g}, \mathbf{y} - \mathbf{x} \rangle + (\rho/2) ||\mathbf{y} - \mathbf{x}||_2^2$ ,  $\forall \mathbf{x}, \mathbf{y} \in \Omega$  and  $\mathbf{g} \in \partial h(\mathbf{x})$ .

**Assumption 2.1** (Strong convexity). f+g has a unique minimizer  $\widehat{\boldsymbol{\theta}} \in \mathbb{R}^p$ , and is  $\rho$ -strongly convex in  $B(\widehat{\boldsymbol{\theta}}, R)$  for some R > 0 and  $\rho > 0$ .

**Assumption 2.2** (Homogeneity).  $\|\nabla^2 f_k(\boldsymbol{\theta}) - \nabla^2 f(\boldsymbol{\theta})\|_2 \leq \delta$  holds for all  $k \in [m]$  and  $\boldsymbol{\theta} \in B(\widehat{\boldsymbol{\theta}}, R)$ .

We will refer to  $\delta$  as a homogeneity parameter. Based on both assumptions, we define

$$\rho_0 = \sup \left\{ c \in [0, \rho] : \{ f_k + g \}_{k=1}^m \text{ are } c\text{-strongly convex in } B(\widehat{\boldsymbol{\theta}}, R) \right\}.$$
 (2.3)

A simple but useful fact is  $\max\{\rho - \delta, 0\} \leq \rho_0 \leq \rho$ . In most interesting problems, the population risk F is smooth and strongly convex on any compact set. When  $\{\mathbf{X}_i\}_{i=1}^N$  are i.i.d. and the total sample size N is large, the empirical risk f concentrates around its population counterpart and inherits nice properties from the latter, making Assumption 2.1 hold easily.

Since  $\{f_k\}_{k=1}^m$  are i.i.d. stochastic approximations of F, they should not be too far away from their average f provided that n is not too small. Assumption 2.2 is a natural way of characterizing this similarity. It is a generalization of the concept " $\delta$ -related functions" for quadratic losses in Arjevani and Shamir (2015). With high probability, it holds with reasonably small  $\delta$  and large R under general conditions (Mei et al., 2018). Large n implies small homogeneity parameter  $\delta$  and thus similar  $\{f_k\}_{k=1}^m$ . Assumption 2.2 always holds with  $\delta = \max_{k \in [m]} \sup_{\theta \in B(\widehat{\theta}, R)} \|\nabla^2 f_k(\theta)\|_2 + \sup_{\theta \in B(\widehat{\theta}, R)} \|\nabla^2 f(\theta)\|_2$ .

In this section, we restrict ourselves to the large-sample regime where the local sample size n is sufficiently large such that  $\rho_0 > \delta \ge 0$ , where  $\rho_0$  is the strong convexity parameter in (2.3). General cases will be discussed in Section 3 where additional local regularization is needed.

The following additional assumption on smoothness of the Hessian matrix of f + g is not necessary for contraction, but it helps us obtain a much stronger result on the contraction rate of Algorithm 2, justifying the power of the simple averaging step.

**Assumption 2.3** (Smoothness of Hessian).  $g \in C^2(\mathbb{R}^p)$ , and there exists  $M \geq 0$  such that  $\|[\nabla^2 f(\boldsymbol{\theta}') + \nabla^2 g(\boldsymbol{\theta}')] - [\nabla^2 f(\boldsymbol{\theta}'') + \nabla^2 g(\boldsymbol{\theta}'')]\|_2 \leq M\|\boldsymbol{\theta}' - \boldsymbol{\theta}''\|_2, \ \forall \boldsymbol{\theta}', \boldsymbol{\theta}'' \in B(\widehat{\boldsymbol{\theta}}, R).$ 

Now for  $k \in [m]$ , define

$$\varphi_k(\boldsymbol{\xi}) = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^p} \left\{ f_k(\boldsymbol{\theta}) + g(\boldsymbol{\theta}) - \langle \nabla f_k(\boldsymbol{\xi}) - \nabla f(\boldsymbol{\xi}), \boldsymbol{\theta} \rangle \right\}.$$

In Algorithm 1, we have  $\boldsymbol{\theta}_{t+1} = \varphi_1(\boldsymbol{\theta}_t)$ ; in Algorithm 2, we have  $\boldsymbol{\theta}_{t+1} = \frac{1}{m} \sum_{k=1}^m \varphi_k(\boldsymbol{\theta}_t)$ . Note that  $\varphi_k(\widehat{\boldsymbol{\theta}}) = \widehat{\boldsymbol{\theta}}$  and  $\widehat{\boldsymbol{\theta}}$  is a fixed point of  $\varphi_k$ . This is the key to success of the algorithms. The following theorem describes the contraction of optimization errors of Algorithms 1 and 2. It is deterministic and non-asymptotic by nature.

**Theorem 2.1.** Let Assumptions 2.1 and 2.2 hold, and  $\rho_0 > \delta \geq 0$ . Consider the iterates  $\{\boldsymbol{\theta}_t\}_{t=0}^{\infty}$  produced by Algorithm 1 or 2, with  $\boldsymbol{\theta}_0 \in B(\widehat{\boldsymbol{\theta}}, R)$ . Then

$$\|\boldsymbol{\theta}_{t+1} - \widehat{\boldsymbol{\theta}}\|_2 \le (\delta/\rho_0) \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2, \quad \forall t \ge 0.$$

In addition, if Assumption 2.3 also holds, then for Algorithm 2 we have

$$\|\boldsymbol{\theta}_{t+1} - \widehat{\boldsymbol{\theta}}\|_{2} \leq \frac{\delta}{\rho_{0}} \|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2} \cdot \min \left\{ 1, \frac{\delta}{\rho} \left( 1 + \frac{M}{\rho_{0}} \|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2} \right) \right\}, \qquad \forall t \geq 0.$$

Theorem 2.1 shows the Q-linear convergence\* of the sequence  $\{\boldsymbol{\theta}_t\}_{t=0}^{\infty}$  generated by both Algorithms 1 and 2. The contraction rate depends explicitly on homogeneity parameter  $\delta$ . With an additional standard assumption on Hessian smoothness, we further show that the averaging step alone in Algorithm 2 is almost as powerful as an optimization step in terms of contraction: The contracting constant will eventually be  $\frac{\delta}{\rho_0} \frac{\delta}{\rho}$ . With negligible computational cost, averaging significantly improves upon individual solutions  $\{\boldsymbol{\theta}_{t,k}\}_{k=1}^m$  by doubling the speed of convergence. In short, one iteration in Algorithm 2 is approximately the same as two iterations in Algorithm 1 under suitable conditions.

The first part of result is a refinement of that in Jordan et al. (2018). In particular, we allow initial estimator to be inaccurate and we have more explicit rates of contraction of optimization errors. This will be demonstrated in Section 2.4 below. The second part points out benefits of the averaging step, which is a novel result.

# 2.4 Multi-step estimators and statistical analysis

While we present the CSL methods as two algorithms, they are really T-step estimators, starting from the initial estimator  $\boldsymbol{\theta}_0$ . The question is then the effect of iterations in the multiple step estimators and the role of the initial estimator. In this section, we show that each iteration makes  $\boldsymbol{\theta}_t$  is closer to the global minimum  $\hat{\boldsymbol{\theta}}$  by a factor of order  $\sqrt{p/n}$  for CSL and by a factor of p/n by the GEL method. This is done by explicitly finding the rate of convergence of  $\delta$ . Thus, through finite steps, the optimization errors are eventually negligible in comparison with statistical errors (assuming N is of order  $(n/p)^a$  for a finite a in typical applications) and the distributed multi-step estimator will work as well as the global minimum as if the data were aggregated in the central server.

<sup>\*</sup>According to Nocedal and Wright (2006), a sequence  $\{\mathbf{x}_n\}_{n=1}^{\infty}$  in  $\mathbb{R}^p$  is said to converge Q-linearly to  $\mathbf{x}^* \in \mathbb{R}^p$  if there exists  $r \in (0,1)$  such that  $\|\mathbf{x}_{n+1} - \mathbf{x}^*\|_2 \le r \|\mathbf{x}_n - \mathbf{x}^*\|_2$  for n sufficiently large.

The deterministic analysis above applies to a wide range of statistical models. As an illustration, we consider the generalized linear model with canonical link, where our samples are i.i.d. pairs  $\{\mathbf{X}_i = (\mathbf{x}_i^T, y_i)^T\}_{i=1}^N$  of covariates and responses and the conditional density of  $y_i$  given  $\mathbf{x}_i$  is given by

$$h(y_i; \mathbf{x}_i, \boldsymbol{\theta}^*) = c(\mathbf{x}_i, y_i) \exp\left(\mathbf{y}_i(\mathbf{x}_i^{\top} \boldsymbol{\theta}^*) - b(\mathbf{x}_i^{\top} \boldsymbol{\theta}^*)\right).$$

Here for simplicity we let the dispersion parameter to be 1 as we do not consider the issue of over-dispersion;  $b(\cdot)$  is some known convex function, and c is a known function such that h is a valid probability density function. The negative log-likelihood of the whole data is an affine transformation of  $f(\boldsymbol{\theta}) = \frac{1}{m} \sum_{k=1}^{m} f_k(\boldsymbol{\theta})$  with

$$f_k(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i \in \mathcal{I}_k} \left[ b(\mathbf{x}_i^{\top} \boldsymbol{\theta}) - y_i(\mathbf{x}_i^{\top} \boldsymbol{\theta}) \right].$$

It's easy to verify that

$$\nabla f_k(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i \in \mathcal{I}_k} [b'(\mathbf{x}_i^{\top} \boldsymbol{\theta}) - y_i] \mathbf{x}_i \quad \text{and} \quad \nabla^2 f_k(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i \in \mathcal{I}_k} b''(\mathbf{x}_i^{\top} \boldsymbol{\theta}) \mathbf{x}_i \mathbf{x}_i^{\top}.$$

Assume that  $\mathbf{x}_i = (1, \mathbf{u}_i^\top)^\top \in \mathbb{R}^p$ , where  $\{\mathbf{u}_i\}_{i=1}^N \subseteq \mathbb{R}^{p-1}$  are i.i.d. random covariate vectors with zero mean and covariance matrix  $\Sigma$ . Suppose there exist universal positive constants  $A_1$ ,  $A_2$  and  $A_3$  such that  $A_1 \leq \|\Sigma\|_2 \leq A_2 p^{A_3}$ . Let  $\Sigma^* = \mathbb{E}(\mathbf{x}_i \mathbf{x}_i^\top) = \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \Sigma \end{pmatrix}$ , g be a deterministic penalty function, and  $F(\theta) = \mathbb{E}f(\theta)$  be the population risk function. Below we impose some standard regularity assumptions.

**Assumption 2.4.** •  $\{\Sigma^{-1/2}\mathbf{u}_i\}_{i=1}^N$  are i.i.d. sub-Gaussian random vectors.

- For all  $x \in \mathbb{R}$ , |b''(x)| and |b'''(x)| are bounded by some constant.
- $\|\boldsymbol{\theta}^*\|_2$  is bounded by some constant.

As in Assumptions 2.1 and 2.3, the following general assumptions is also needed for our analysis. Here R is some positive quantity that satisfies  $R < A_4 p^{A_5}$  for some universal constants  $A_4$  and  $A_5$ .

**Assumption 2.5.** There exists a universal constant  $\rho > 0$  such that (F + g) is  $\rho$ -strongly convex in  $B(\theta^*, 2R)$ .

The following smoothness assumption is only needed for part of our theory; it is used to show that the averaging step in Algorithm 2 can significantly enhance the accuracy.

**Assumption 2.6.**  $g \in C^2(\mathbb{R}^p)$ , and there exists a universal constant  $M \geq 0$  such that  $\|[\nabla^2 F(\boldsymbol{\theta}') + \nabla^2 g(\boldsymbol{\theta}')] - [\nabla^2 F(\boldsymbol{\theta}'') + \nabla^2 g(\boldsymbol{\theta}'')]\|_2 \leq M\|\boldsymbol{\theta}' - \boldsymbol{\theta}''\|_2, \forall \boldsymbol{\theta}', \boldsymbol{\theta}'' \in B(\boldsymbol{\theta}^*, 2R).$ 

Under the model assumptions above, we can explicitly determine rate for  $\delta$  in Assumption 2.2. In particular, we will show in Section A.6 that

$$\max_{k \in [m]} \max_{\boldsymbol{\theta} \in B(\widehat{\boldsymbol{\theta}}, R)} \|\nabla^2 f_k(\boldsymbol{\theta}) - \nabla^2 f(\boldsymbol{\theta})\|_2 = O_{\mathbb{P}} \left( \|\boldsymbol{\Sigma}\|_2 \sqrt{\frac{p(\log p + \log N)}{n}} \right),$$

provided that  $n \geq cp$  for an arbitrary positive constant c. Therefore, with high probability,  $\delta \approx \|\mathbf{\Sigma}\|_2 \sqrt{p(\log p + \log N)/n}$ . Omitting the logarithmic terms, we see that the contraction factor is approximately  $\kappa \sqrt{p/n}$ , where  $\kappa \triangleq \|\mathbf{\Sigma}\|_2/\rho$  can be viewed as condition number. This rate is more explicit on p and  $\kappa$  than that in Jordan et al. (2018), where finite p and  $\kappa$  are assumed. In addition, with a smooth regularization, Algorithm 2 benefits from the averaging step in that it improves the contraction rate to approximately  $\kappa^2 p/n$ .

Let  $\theta_t$  be the t-th step estimator of Algorithm 1 or Algorithm 2 with some initialization  $\theta_0 \in B(\widehat{\theta}, R)$ . It is clear that the statistical error of the estimator  $\theta_t$  is upper bounded by its optimization error and the statistical error of  $\widehat{\theta}$ :

$$\|oldsymbol{ heta}_t - oldsymbol{ heta}^*\|_2 \leq \|oldsymbol{ heta}_t - \widehat{oldsymbol{ heta}}\|_2 + \|\widehat{oldsymbol{ heta}} - oldsymbol{ heta}^*\|_2.$$

The second term is well-studied in statistics, which is of order  $O_{\mathbb{P}}(\sqrt{p/N})$  under mild conditions. The following theorem controls the magnitude of the first term, which is the optimization error.

**Theorem 2.2.** Suppose that Assumptions 2.4 and 2.5 hold and with probability tending to one  $\theta_0 \in B(\widehat{\theta}, R)$  for some  $R > \|\widehat{\theta} - \theta^*\|_2$ . For Algorithms 1 and 2, we have

$$\|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2 = O_{\mathbb{P}}(\eta^{t/2} \|\boldsymbol{\theta}_0 - \widehat{\boldsymbol{\theta}}\|_2), \quad \forall t \ge 0,$$

where  $\eta = \kappa^2 p(\log N)/n$ . In addition, let Assumption 2.6 also hold. There exists some constant C such that for Algorithm 2 we have

$$\|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2 = O_{\mathbb{P}}(\eta^{t-t_0} \|\boldsymbol{\theta}_{t_0} - \widehat{\boldsymbol{\theta}}\|_2), \qquad \forall t \ge t_0, \tag{2.4}$$

where 
$$t_0 = \lceil \frac{2 \log(CMR/\rho)}{\log(1/\eta)} \rceil$$
.

Theorem 2.2 explicitly describes how Algorithms 1 and 2 depend on structural parameters of the problem. First, the condition  $\boldsymbol{\theta}_0 \in B(\widehat{\boldsymbol{\theta}}, R)$  on initialization is mild, since  $\widehat{\boldsymbol{\theta}}$  is usually a consistent estimate and  $\|\boldsymbol{\theta}^*\|_2$  is bounded (Assumption 2.4). In contract with Jordan et al. (2018), we allow inaccurate initial value such as  $\boldsymbol{\theta}_0 = \mathbf{0}$  and give more explicit rates of contraction even when p and  $\kappa$  diverge.

In contrast to fixed contraction derived by Shamir et al. (2014), Theorem 2.2 explains the significant benefits of large local sample size even in the presence of a non-smooth penalty: optimization error is shrunken by a factor that converges to zero explicitly in n. When p and  $\kappa$  are bounded,  $\eta = O((\log N)/n)$  and within a finite step t, the optimization error  $\eta^{t/2}$  can be much smaller than statistical error  $\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_2 = O_{\mathbb{P}}(\sqrt{p/N})$ , so long as  $n \geq N^b$  for some  $0 < b \leq 1$ .

The accuracy of initial estimator  $\boldsymbol{\theta}_0$  helps reducing the number of iterations. As an example of this, consider the simple average of individual estimator for node machine as  $\boldsymbol{\theta}_0$  for smooth loss function with no regularization. By Corollary 2 in Zhang et al. (2013), this simple divide-and-conquer estimator has accuracy  $\|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*\|_2 = O_{\mathbb{P}}(\max\{\sqrt{\frac{p}{N}}, \frac{p}{n}, \frac{\kappa\sqrt{p\log p}}{n}\})$ . Using the explicit expression of  $\eta$ , we can easily show that the one-step estimator  $\boldsymbol{\theta}_1$  obtained by Algorithm 1 behaves the same as the global minimizer  $\hat{\boldsymbol{\theta}}$  if the local sample size

$$n^3 \gg N\kappa^2 p \log N(p + \kappa^2 \log p). \tag{2.5}$$

In this case, the local optimization in Algorithm 1 can further be replaced by using the explicit one-step estimator as in Bickel (1975) and Jordan et al. (2018). A similar remark applies to the GEL method (Algorithm 2).

# 3 Distributed estimation in general regimes

Algorithms 1 and 2 and their analysis in the last section are built upon the large-sample regime, with sufficiently strong convexity of  $\{f_k + g\}_{k=1}^m$  and small discrepancy between them. This requires the local sample size n to be large enough, which may not be the case in practice. Even worse, the required local sample size depends on structural parameters, making such a condition unverifiable. Our numerical experiments confirm the instability of Algorithms 1 and 2 even for moderate n. A naive method of remedy is to add strict convex quadratic regularization  $q(\theta)$ . While this remedy can make the algorithm converges rapidly,

the nonadaptive nature of  $q(\theta)$  will make the convergence to a wrong target. Instead of using a fixed q, we will adjust the regularization function according to current solutions. The idea stems from the proximal point algorithm (Rockafellar, 1976; Parikh and Boyd, 2014).

## 3.1 Distributed approximate proximal point algorithms

**Definition 3.1.** For any convex function  $h : \mathbb{R}^p \to \mathbb{R}$ , define the proximal mapping  $\operatorname{prox}_h : \mathbb{R}^p \to \mathbb{R}^p$ ,  $\mathbf{x} \mapsto \operatorname{argmin}_{\mathbf{y} \in \mathbb{R}^p} \{ h(\mathbf{y}) + \|\mathbf{y} - \mathbf{x}\|_2^2 / 2 \}$ .

For a given  $\alpha > 0$ , the proximal point algorithm for minimizing h iteratively computes

$$\mathbf{x}_{t+1} = \operatorname{prox}_{\alpha^{-1}h}(\mathbf{x}_t) = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^p} \{ h(\mathbf{x}) + (\alpha/2) \|\mathbf{x} - \mathbf{x}_t\|_2^2 \}, \quad \forall t \ge 0,$$

starting from some initial value  $\mathbf{x}_0$ . Under mild conditions,  $\{\mathbf{x}_t\}_{t=0}^{\infty}$  converges linearly to some  $\hat{\mathbf{x}} \in \operatorname{argmin}_{\mathbb{R}^p} h(\mathbf{x})$  (Rockafellar, 1976).

Now we take h = f + g and write the proximal point iteration for our problem (2.2):

$$\boldsymbol{\theta}_{t+1} = \operatorname{prox}_{\alpha^{-1}(f+g)}(\boldsymbol{\theta}_t) = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^p} \left\{ f(\boldsymbol{\theta}) + g(\boldsymbol{\theta}) + \frac{\alpha}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|_2^2 \right\}.$$
 (3.1)

Each iteration (3.1) is a distributed optimization problem, whose object function is not available to node machines. But it can be solved by Algorithms 1 and 2. Specifically, suppose we have already obtained  $\boldsymbol{\theta}_t$  and aim for  $\boldsymbol{\theta}_{t+1}$  in (3.1). Letting  $\tilde{g}(\boldsymbol{\theta}) = g(\boldsymbol{\theta}) + (\alpha/2) \|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|_2^2$ , Algorithm 2 starting from  $\tilde{\boldsymbol{\theta}}_0 = \boldsymbol{\theta}_t$  produces iterations over  $s = 0, 1, \cdots$ 

$$\tilde{\boldsymbol{\theta}}_{s,k} = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^p} \left\{ f_k(\boldsymbol{\theta}) + \tilde{g}(\boldsymbol{\theta}) + \langle \nabla f_k(\tilde{\boldsymbol{\theta}}_s) - \nabla f(\tilde{\boldsymbol{\theta}}_s), \boldsymbol{\theta} \rangle \right\}, \qquad k \in [m],$$

$$\tilde{\boldsymbol{\theta}}_{s+1} = \frac{1}{m} \sum_{k=1}^m \tilde{\boldsymbol{\theta}}_{s,k}.$$

When  $\alpha + \rho_0 > \delta$ ,  $\{\tilde{\boldsymbol{\theta}}_s\}_{s=0}^{\infty}$  converges Q-linearly to  $\boldsymbol{\theta}_{t+1}$ . On the other hand, there is no need to solve (3.1) exactly, as  $\operatorname{prox}_{\alpha^{-1}(f+g)}(\boldsymbol{\theta}_t)$  is merely an intermediate quantity for computing  $\hat{\boldsymbol{\theta}}$ . We therefore only run one iteration of the GEL Algorithm 2 and use the resulting approximate solution as  $\boldsymbol{\theta}_{t+1}$ . This simplifies the algorithm, reducing double loops to a single loop, and enhances statistical interpretation of the method as a multi-step estimator. However, it makes technical arguments more challenging. Similarly, we may also use one step of Algorithm 1 to compute the inexact proximal update.

The above discussions lead us to propose two Communication-Efficient Accurate Statistical Estimators (CEASE) in Algorithms 3 and 4, which use the proximal point algorithm as the backbone and obtain inexact updates in a distributed manner. They are regularized versions of Algorithms 1 and 2, with an additional proximal term in the objective functions. The term reduces relative differences of the local loss functions on individual machines, and is particularly crucial for convergence when  $\{f_k\}_{k=1}^m$  are not similar enough. Ideas from the proximal point algorithm have appeared in the literature of distributed stochastic optimization for different purposes such as accelerating first-order algorithms (Lee et al., 2017a) and regularizing sizes of updates (Wang et al., 2017b).

#### Algorithm 3 Communication-Efficient Accurate Statistical Estimators (CEASE)

**Input**: Initial value  $\theta_0$ , regularizer  $\alpha \geq 0$ , number of iterations T. For  $t = 0, 1, 2, \dots, T - 1$ :

- Each machine evaluates  $\nabla f_k(\boldsymbol{\theta}_t)$  and sends to the 1<sup>st</sup> machine;
- The 1<sup>st</sup> machine computes  $\nabla f(\boldsymbol{\theta}_t) = \frac{1}{m} \sum_{k=1}^m \nabla f_k(\boldsymbol{\theta}_t)$  and

$$\boldsymbol{\theta}_{t+1} = \operatorname{argmin}_{\boldsymbol{\theta}} \left\{ f_1(\boldsymbol{\theta}) + g(\boldsymbol{\theta}) - \langle \nabla f_k(\boldsymbol{\theta}_t) - \nabla f(\boldsymbol{\theta}_t), \boldsymbol{\theta} \rangle + \frac{\alpha}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|_2^2 \right\},$$

and broadcasts to other machines.

#### Output: $\theta_T$ .

In each iteration, Algorithm 3 has one round of communication and one optimization problem to solve. Although Algorithm 4 has two rounds of communication per iteration, only one round involves parallel optimization and the other is simply averaging. We will compare their theoretical guarantees as well as practical performances in the sequel.

# 3.2 Contraction of optimization errors

Theorem 3.1 gives contraction guarantees for Algorithms 3 and 4.

**Theorem 3.1.** Let Assumptions 2.1 and 2.2 hold. Consider the multi-step estimators  $\{\boldsymbol{\theta}_t\}_{t=0}^T$  generated by Algorithm 3 or 4. Suppose that  $\boldsymbol{\theta}_0 \in B(\widehat{\boldsymbol{\theta}}, R/2)$  and  $[\delta/(\rho_0 + \alpha)]^2 < \rho/(\rho + 2\alpha)$ .

#### Algorithm 4 CEASE with averaging

**Input**: Initial value  $\theta_0$ , regularizer  $\alpha \geq 0$ , number of iterations T. For  $t = 0, 1, 2, \dots, T - 1$ :

- Each machine evaluates  $\nabla f_k(\boldsymbol{\theta}_t)$  and sends to the central processor;
- The central processor computes  $\nabla f(\boldsymbol{\theta}_t) = \frac{1}{m} \sum_{k=1}^m \nabla f_k(\boldsymbol{\theta}_t)$  and broadcasts to machines;
- Each machine computes

$$\boldsymbol{\theta}_{t,k} = \operatorname{argmin}_{\boldsymbol{\theta}} \left\{ f_k(\boldsymbol{\theta}) + g(\boldsymbol{\theta}) - \langle \nabla f_k(\boldsymbol{\theta}_t) - \nabla f(\boldsymbol{\theta}_t), \boldsymbol{\theta} \rangle + \frac{\alpha}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|_2^2 \right\}$$

and sends to the central processor;

• The central processor computes  $\theta_{t+1} = \frac{1}{m} \sum_{k=1}^{m} \theta_{t,k}$  and broadcasts to machines.

#### Output: $\theta_T$ .

• For both Algorithms 3 and 4, we have

$$\|\boldsymbol{\theta}_{t+1} - \widehat{\boldsymbol{\theta}}\|_{2} \le \|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2} \cdot \frac{\frac{\delta}{\rho_{0} + \alpha} \sqrt{\rho^{2} + 2\alpha\rho} + \alpha}{\rho + \alpha}, \qquad 0 \le t \le T - 1;$$
 (3.2)

• If Assumption 2.3 also holds, then for Algorithm 4 we have

$$\|\boldsymbol{\theta}_{t+1} - \widehat{\boldsymbol{\theta}}\|_{2} \le \|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2} \cdot \frac{\gamma_{t}\sqrt{\rho^{2} + 2\alpha\rho} + \alpha}{\rho + \alpha}, \qquad 0 \le t \le T - 1,$$
 (3.3)

where we define  $\gamma_t = \frac{\delta}{\rho_0 + \alpha} \cdot \min\{1, \frac{\delta}{\rho + \alpha}(1 + \frac{M}{\rho_0 + \alpha} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2)\};$ 

• Both multiplicative factors in (3.2) and (3.3) are strictly less than 1.

In the contraction factor in (3.2), the two summands  $\frac{\delta}{(\rho_0 + \alpha)^2} \sqrt{\rho^2 + 2\alpha\rho}$  and  $\frac{\alpha}{\rho + \alpha}$  come from the error of the inexact proximal update  $\|\boldsymbol{\theta}_{t+1} - \operatorname{prox}_{\alpha^{-1}(f+g)}(\boldsymbol{\theta}_t)\|_2$  and the residual of the proximal point  $\|\operatorname{prox}_{\alpha^{-1}(f+g)}(\boldsymbol{\theta}_t) - \widehat{\boldsymbol{\theta}}\|_2$ , respectively. Similar results hold for (3.3).

Theorem 3.1 justifies the linear convergence of Algorithms 3 and 4 under quite general settings. The local loss functions  $\{f_k\}_{k=1}^m$  just need to be convex and smooth, and the convex penalty g is allowed to be non-smooth, e.g. the  $\ell_1$  norm. On the contrary, most algorithms for distributed statistical estimation are only designed for smooth problems, and many of them are only rigorously studied when the loss functions are quadratic or self-concordant

(Shamir et al., 2014; Zhang and Xiao, 2015; Wang et al., 2017b). This is another important aspect of our contributions.

Note that the convergence of the vanilla Algorithms 1 and 2 hinges on the homogeneity assumption  $\rho_0 > \delta$  in Theorem 2.1, i.e. the functions  $\{f_k\}_{k=1}^m$  must be similar enough. In the statistical setting, this requires n to be large. Algorithms 3 and 4 no longer need such a condition and converge linearly as long as  $[\delta/(\rho_0 + \alpha)]^2 < \rho/(\rho + 2\alpha)$ , which is guaranteed to hold by choosing sufficiently large  $\alpha$ . Hence proper regularization provides a safety net for the algorithms. Corollary 3.1 below gives a guideline for choosing  $\alpha$  to make Algorithms 3 and 4 converge.

Corollary 3.1. Let Assumptions 2.1 and 2.2 hold,  $\boldsymbol{\theta}_0 \in B(\widehat{\boldsymbol{\theta}}, R/2)$ , and  $\{\boldsymbol{\theta}_t\}_{t=0}^T$  be the iterates of Algorithm 3 or 4. With any  $\alpha \geq 4\delta^2/\rho$ , both algorithms converge with contraction factors in (3.2) and (3.3) bounded by  $(1 - \frac{\rho}{10(\alpha+\rho)})$ . Hence to reach the statistically negligible accuracy of  $O((\frac{p}{N})^{1/2+\epsilon_0})$  for a constant  $\epsilon_0 > 0$ , we need at most  $T = O((1 + \frac{\alpha}{\rho})\log(\|\boldsymbol{\theta}_0 - \widehat{\boldsymbol{\theta}}\|_2 \cdot \frac{N}{p}))$  iterations.

Consider again the case where the local loss functions have small relative difference  $\delta/\rho$ . In this case, Theorem 2.1 states that the contraction factors for unregularized versions ( $\alpha = 0$ ) of Algorithms 3 and 4 are in the same order of  $\delta/\rho$ . The following corollary tells us how large  $\alpha$  can be so that the contraction factors are still of that order. It provides an upper bound for the amount of regularization to make the algorithms converge rapidly in nice scenarios.

Corollary 3.2. Let Assumptions 2.1 and 2.2 hold,  $\theta_0 \in B(\widehat{\theta}, R/2)$ , and suppose  $\alpha \leq C\delta^2/\rho$  for some constant C. There exist constants  $C_1$  and  $C_2$  such that the followings hold when  $\delta/\rho$  is sufficiently small:

• Algorithms 3 and 4 have the contraction property

$$\|\boldsymbol{\theta}_{t+1} - \widehat{\boldsymbol{\theta}}\|_2 \le C_1 \delta / \rho \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2, \qquad 0 \le t \le T - 1;$$

• if Assumption 2.3 also holds and  $\|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2 \leq \rho/M$ , then

$$\|\boldsymbol{\theta}_{t+1} - \widehat{\boldsymbol{\theta}}\|_2 \le C_2(\delta/\rho)^2 \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2$$

holds for Algorithm 4.

Consequently,  $T = O(\log(\|\boldsymbol{\theta}_0 - \widehat{\boldsymbol{\theta}}\|_2 \cdot \frac{N}{p})/\log(\frac{\rho}{\delta}))$  suffices for both algorithms to achieve a statistically negligible accuracy of  $O((\frac{p}{N})^{1/2+\epsilon_0})$  with a small constant  $\epsilon_0 > 0$ .

Corollary 3.2 suggests choosing a small regularizer  $\alpha = O(\delta^2/\rho)$  when  $\delta/\rho$  is small. The contraction factors in Corollary 3.2 go to zero if  $\delta/\rho$  does, indicating both algorithms' ability to utilize the similarity among local loss functions. With a regularizer  $\alpha$  up to the order of  $\delta^2/\rho$ , the contraction factors are essentially the same as those of the unregularized ( $\alpha = 0$ ) algorithms. If (f+g) is smooth and  $\theta_t$  is reasonably close to  $\hat{\theta}$ , then Corollary 3.2 shows that each iteration of Algorithm 4 is roughly equivalent to two iterations of Algorithm 3, although the former only has one round of optimization. The averaging step in Algorithm 4 reduces the error as much as the optimization step, while taking much less time. In this case, Algorithm 4 is preferable, and our numerical experiments also confirm this.

By combining Corollaries 3.1 and 3.2, we get

$$\alpha \simeq \delta^2/\rho$$

as a default choice for Algorithms 3 and 4 to become fast and robust. They are reliable in general cases and efficient in nice cases.

According to the results above, Algorithms 3 and 4 achieve  $\varepsilon$ -accuracy within

$$O\left(\max\{1, (\delta/\rho)^2\} \log\left(\frac{\|\boldsymbol{\theta}_0 - \widehat{\boldsymbol{\theta}}\|_2}{\varepsilon}\right)\right)$$
 (3.4)

rounds of communication. In contrast, the distributed accelerated gradient descent requires  $O(\sqrt{\kappa_0}\log(\frac{\|\theta_0-\widehat{\theta}\|_2}{\varepsilon}))$  rounds of communication to achieve  $\varepsilon$ -accuracy (Shamir et al., 2014), with  $\kappa_0$  being the condition number of (f+g), which does not take advantage of sample size n. As long as  $\delta/\rho \ll \kappa_0^{1/4}$ , our Algorithms 3 and 4 communicate less than the distributed accelerated gradient descent. This is achieved by leveraging the similarity among  $\{f_k\}_{k=1}^m$ . And again, our general results for Algorithms 3 and 4 also apply to the case with nonsmooth penalty functions while those for distributed accelerated gradient descent do not.

For unregularized empirical risk minimization, i.e. g = 0 in (2.2), Algorithm 4 reduces to an extension or a useful case of the DANE algorithm (Shamir et al., 2014). DANE is motivated from the mirror descent and only deals with smooth optimization problems. Theoretical analysis of DANE beyond quadratic loss require extremal choice of tuning parameters and does not show any advantage over distributed implementation of the gradient descent. On the other hand, we derive Algorithm 4 from the proximal point algorithm, handling both smooth and nonsmooth problems. This new perspective leads to sharp analysis of Algorithm 4 along with suggestions on choosing the tuning parameter  $\alpha$ . As a by-product,

we close a gap in the theory of DANE in non-quadratic settings. Our analysis techniques are potentially useful for other distributed optimization algorithms, especially when the loss is not quadratic.

## 3.3 Multi-step estimators and their statistical properties

#### 3.3.1 General case

Consider again the generalized linear model with the canonical link as in Section 2.4. In the following theorem, we specify the correct order of the regularization parameter  $\alpha$  such that Algorithms 3 and 4 not only overcome the difficulties with a small local sample size n, but also inherit all the advantages of previous algorithms in the large-n regime. This leverages upon the explicit homogeneity rate  $\delta$ .

**Theorem 3.2.** Suppose that Assumptions 2.4 and 2.5 hold, and with high probability  $\boldsymbol{\theta}_0 \in B(\widehat{\boldsymbol{\theta}}, R/2)$ . Let  $\eta = \kappa^2(\log N)p/n$  and  $\kappa = \|\boldsymbol{\Sigma}\|_2/\rho$ . For any  $c_1, c_2 > 0$ , there exists C > 0 such that the followings hold with high probability:

• if  $n \ge c_1 p$  and  $\alpha \ge C \rho \eta$ , then both Algorithms 3 and 4 have linear convergence

$$\|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2 \le \left[1 - \frac{\rho}{10(\alpha + \rho)}\right]^t \|\boldsymbol{\theta}_0 - \widehat{\boldsymbol{\theta}}\|_2, \quad \forall t \ge 0;$$

• if  $\eta$  is sufficiently small and  $\alpha \leq c_2 \rho \eta$ , then for both algorithms

$$\|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2 = O_{\mathbb{P}}(\eta^{t/2} \|\boldsymbol{\theta}_0 - \widehat{\boldsymbol{\theta}}\|_2), \quad \forall t \ge 0;$$

in addition, if Assumption 2.6 also holds, then for Algorithm 4 we have

$$\|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2 = O_{\mathbb{P}}(\eta^{t-t_0} \|\boldsymbol{\theta}_{t_0} - \widehat{\boldsymbol{\theta}}\|_2), \quad \forall t \ge t_0,$$

where 
$$t_0 = \lceil \frac{2 \log(CMR/\rho)}{\log(1/\eta)} \rceil$$
.

For many big-data problems of interest, it is reasonable to assume that n/p is bounded away from 0 by some small constant. Then Theorem 3.2 indicates that by choosing  $\alpha \simeq \rho \eta$ , Algorithms 3 and 4 inherit all the merits of Algorithms 1 and 2 in the large n regime – fast linear contraction of rate  $\sqrt{\eta} = \kappa \sqrt{p(\log N)/n}$ , and for Algorithm 4, a even faster rate of  $\eta = \kappa^2 p(\log N)/n$  to  $\widehat{\boldsymbol{\theta}}$  when the loss functions and the penalty are smooth. These facts

also guarantee that Algorithms 3 and 4 reach the statistical efficiency in  $O(\frac{\log(\|\boldsymbol{\theta}_0 - \widehat{\boldsymbol{\theta}}\|_2 \sqrt{N/p})}{\log(1/\eta)})$  iterations. While it is hard to check whether n is sufficiently large in practice, proper choice of  $\alpha$  always guarantees linear convergence, and the contraction rates adapt to the sample size n. In this way, Algorithms 3 and 4 perfectly resolve the main issue of their vanilla versions.

Similar to the discussion at the end of Seciton 2.4,  $\theta_t$  produces by CEASE and CEASE with averaging can be regarded as multi-step statistical estimators. As the contraction of optimization error is at least at the order of  $\sqrt{\eta}$ , it only takes finite steps to achieve negligible optimization error and hence achieves statistical efficiency in typical statistical applications. In particular, when n satisfies (2.5), the one-step estimator from the one-shot average (Zhang et al., 2013) is statistically efficient and its asymptotic inference follows from that based on the empirical minimizer  $\hat{\theta}$  based on all the data.

#### 3.3.2 Quadratic loss

We can get stronger results in the specific case of distributed linear regression: the contraction rate has nearly no dependence on the conditional number  $\kappa$ . This is demonstrated by leveraging on the analytic solutions. In this case, the  $k^{th}$  machine defines a quadratic loss function

$$\frac{1}{2n} \sum_{i \in \mathcal{I}_k} (y_i - \mathbf{x}_i^{\top} \boldsymbol{\theta})^2 = \frac{1}{2} \boldsymbol{\theta}^{\top} \widehat{\boldsymbol{\Sigma}}_k \boldsymbol{\theta} - \widehat{\mathbf{w}}_k^{\top} \boldsymbol{\theta} + \frac{1}{2n} \sum_{i \in \mathcal{I}_k} y_i^2,$$

where  $\widehat{\boldsymbol{\Sigma}}_k = \frac{1}{n} \sum_{i \in \mathcal{I}_k} \mathbf{x}_i \mathbf{x}_i^{\top}$  and  $\widehat{\mathbf{w}}_k = \frac{1}{n} \sum_{i \in \mathcal{I}_k} \mathbf{x}_i y_i$ . Let  $f(\boldsymbol{\theta}) = \frac{1}{2} \boldsymbol{\theta}^{\top} \widehat{\boldsymbol{\Sigma}} \boldsymbol{\theta} - \widehat{\mathbf{w}}^{\top} \boldsymbol{\theta}$ . Without loss of geneality we write  $\mathbf{x}_i = (1, \mathbf{u}_i^{\top})^{\top} \in \mathbb{R}^p$ .

Assumption 3.1. •  $\mathbb{E}\mathbf{u}_i = 0$  and  $\mathbb{E}(\mathbf{u}_i\mathbf{u}_i^{\top}) = \mathbf{\Sigma} \succ 0$ .  $\{\mathbf{\Sigma}^{-1/2}\mathbf{u}_i\}_{i=1}^N$  are i.i.d. sub-Gaussian random vectors with bounded  $\|\mathbf{\Sigma}^{-1/2}\mathbf{u}_i\|_{\psi_2}$ .

- The minimum eigenvalue  $\lambda_{\min}(\Sigma)$  is bounded away from zero.
- $N/\text{Tr}(\Sigma) \ge C > 0$  and  $n/\log m \ge c > 0$  where C and c are constants.

For the least-squares, Algorithm 3 admits a close-form:

$$\boldsymbol{\theta}_{t+1} = [\mathbf{I} - (\widehat{\boldsymbol{\Sigma}}_1 + \alpha \mathbf{I})^{-1} \widehat{\boldsymbol{\Sigma}}] \boldsymbol{\theta}_t + (\widehat{\boldsymbol{\Sigma}}_1 + \alpha \mathbf{I})^{-1} \widehat{\mathbf{w}},$$

and so does Algorithm 4:

$$\boldsymbol{\theta}_{t+1,k} = [\mathbf{I} - (\widehat{\boldsymbol{\Sigma}}_k + \alpha \mathbf{I})^{-1} \widehat{\boldsymbol{\Sigma}}] \boldsymbol{\theta}_t + (\widehat{\boldsymbol{\Sigma}}_k + \alpha \mathbf{I})^{-1} \widehat{\mathbf{w}},$$

$$\boldsymbol{\theta}_{t+1} = \left(\mathbf{I} - \frac{1}{m} \sum_{k=1}^{m} (\widehat{\boldsymbol{\Sigma}}_k + \alpha \mathbf{I})^{-1} \widehat{\boldsymbol{\Sigma}}\right) \boldsymbol{\theta}_t + \frac{1}{m} \sum_{k=1}^{m} (\widehat{\boldsymbol{\Sigma}}_k + \alpha \mathbf{I})^{-1} \widehat{\mathbf{w}}.$$

Intuitively, the averaging step in Algorithm 4 reduces variance and accelerates convergence. Below we study Algorithm 4 with the help of these analytical expressions. In the large sample regime, we achieve a contraction factor of O(p/n) without any condition number; in the general regime, linear convergence is still guaranteed.

**Theorem 3.3.** Suppose Assumption 3.1 holds and n/p is bounded away from zero. Then, there exist positive constants  $C_1, C_2$  and  $C_3$  such that when (i)  $n \geq C_1p$  and  $\alpha \geq 0$  or (ii)  $\alpha \geq C_1 \text{Tr}(\Sigma)/n$ , with probability tending to 1,

$$\|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2 \le 2\sqrt{\kappa} \ \eta^t \|\boldsymbol{\theta}_0 - \widehat{\boldsymbol{\theta}}\|_2, \qquad \forall t \ge 0,$$
 (3.5)

where 
$$\kappa = \lambda_{\max}(\Sigma)/\lambda_{\min}(\Sigma)$$
 and  $\eta = 1 - \frac{1 - \min\{1/2, C_2 p/n\}}{1 + C_3 \alpha}$ .

Theorem 3.3 reveals the following remarkable facts about Algorithm 4: No matter what relationship n and p have, proper regularization always guarantees linear convergence, and the rate exhibits a smooth transition as p/n grows. Hence we can handle the distributed statistical estimation problem without assuming large enough n, overcoming the difficulty of other algorithms in literature (Zhang et al., 2013; Battey et al., 2015; Jordan et al., 2018).

If n/p is large enough, the regularization is not necessary, but choosing  $\alpha \approx p/n$  does not hurt much. This is because we can control the contraction factor as:

$$1 - \frac{1 - C_2 p/n}{1 + C_3 \alpha} = \frac{C_3 \alpha + C_2 p/n}{1 + C_3 \alpha} = O(p/n).$$

When n/p is not that large, most distributed statistical estimation procedures fail. By choosing  $\alpha = \tilde{C} \operatorname{Tr}(\Sigma)/n$  for  $\tilde{C} > C_1$  (see Condition (ii) of Theorem 3.3) we still have linear convergence with contraction factor at most

$$1 - \frac{1 - 1/2}{1 + C_3 \alpha} = 1 - \frac{1}{2 + 2C_3 \tilde{C} \text{Tr}(\mathbf{\Sigma})/n} < 1.$$

In most situations of interest we have  $\text{Tr}(\Sigma) \simeq p$  (even for pervasive factor models). Therefore we see that  $\alpha \simeq p/n$  is a universal and adaptive choice of regularization over all the possible relation between n and p.

Another benefit of the Algorithms is that the condition number  $\kappa$  has only logarithmic

effect on the iteration complexity, and the contraction factor in Theorem 3.3 does not depend on  $\kappa$  at all. This is in stark contrast to the analysis under the same setting in Shamir et al. (2014), and helps relax the commonly used boundedness assumption on the condition number in Zhang et al. (2013), Battey et al. (2015), Jordan et al. (2018), among others. It is worth mentioning that Wang et al. (2018a) derive similar results for distributed linear regression when the local sample size n is sufficiently large.

# 4 Numerical experiments

## 4.1 Synthetic data

We first conduct distributed logistic regression to illustrate the effect of local sample size and initialization on the convergence. We keep the total sample size N=10000 and the dimensionality p=100 fixed, and generate the i.i.d. data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  as follows:  $\mathbf{x}_i = (1, \mathbf{u}_i^{\mathsf{T}})^{\mathsf{T}}$  with  $\mathbf{u}_i \sim N(\mathbf{0}_p, \mathbf{\Sigma})$  and  $\mathbf{\Sigma} = \mathrm{diag}(10, 5, 2, 1 \cdots 1) \in \mathbb{R}^{p \times p}$ ;  $\mathbb{P}(y_i = 1) = 1 - \mathbb{P}(y_i = 0) = \frac{e^{\mathbf{x}_i^{\mathsf{T}} \boldsymbol{\theta}^*}}{1+e^{\mathbf{x}_i^{\mathsf{T}} \boldsymbol{\theta}^*}}$  where  $\boldsymbol{\theta}^* \in \mathbb{R}^{p+1}$  is a random vector with norm 3 whose direction is chosen uniformly at random from the sphere. We use the natural logarithm of the estimation error  $\|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|_2$  to measure the performance of different algorithms, including multiple versions of the CEASE algorithms, GIANT (Wang et al., 2018a), ADMM (Boyd et al., 2011) and accelerated gradient descent (Nesterov, 1983).

Figure 1 shows how the estimation errors evolve with iterations, and the numerical results are the average values over 100 independent runs. The regimes "large n", "moderate n" and "small n" refer to (n, m) = (2000, 5), (1000, 10) and (250, 40); "zero initialization" and "good initialization" refer to  $\boldsymbol{\theta}_0 = \mathbf{0}$  (bottom panel) and  $\bar{\boldsymbol{\theta}}$  (top panel), respectively. Here  $\bar{\boldsymbol{\theta}}$  is the one-shot distributed estimator (Zhang et al., 2013) that averages the individual estimators on local machines.

With proper regularization, the two CEASE algorithms are the only ones that converge rapidly in all scenarios. The purely deterministic methods ADMM (Boyd et al., 2011) and accelerated gradient descent (Nesterov, 1983) are also reliable but slow. Other distributed algorithms like unregularized CEASE and GIANT (Wang et al., 2018a) easily fail when the local sample size is small or the initialization is uninformative. In addition, the CEASE with averaging (Algorithm 4) is superior to the one without averaging (Algorithm 3). For example, when (n, m) = (1000, 10), the averaged CEASE with  $\alpha = 0$  converges while the one without averaging does not. Hence the averaging step leads to better performance.

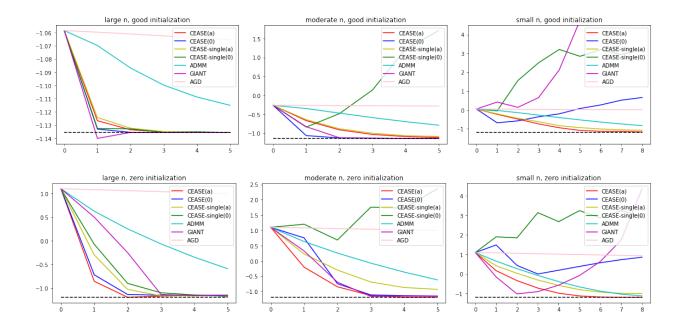


Figure 1: Impacts of local sample size and initialization on convergence. The x-axis and y-axis are the number of iterations and  $\log \|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|_2$ . The dashed lines show the error of the minimizer of the overall loss function. Top panel uses  $\bar{\boldsymbol{\theta}}$  as the initial value and bottom uses 0 as the initial value. CEASE(a) and CEASE(0) refer to Algorithm 4 with  $\alpha = 0.15p/n$  and 0; CEASE-single(a) and CEASE-single(0) refer to Algorithm 3 with  $\alpha = 0.15p/n$  and 0, respectively.

Next, we use  $\ell_1$ -regularized logistic regression to validate the efficacy of our algorithms in the presence of a nonsmooth penalty. We fix the total sample size N=5000 and the dimensionality p=1000, and generate the i.i.d. data  $\{(\mathbf{x}_i,y_i)\}_{i=1}^N$  as follows:  $\mathbf{x}_i=(1,\mathbf{u}_i^\top)^\top$  with  $\mathbf{u}_i \sim N(\mathbf{0}_p,\mathbf{I}_p)$ ;  $\mathbb{P}(y_i=1)=1-\mathbb{P}(y_i=0)=\frac{e^{\mathbf{x}_i^\top\boldsymbol{\theta}^*}}{1+e^{\mathbf{x}_i^\top\boldsymbol{\theta}^*}}$  where  $\boldsymbol{\theta}^*=(\mathbf{1}_{10}^\top,\mathbf{0}_{991}^\top)^\top/\sqrt{2}\in\mathbb{R}^{p+1}$ . We define the penalty function  $g(\boldsymbol{\theta})=\lambda\|\boldsymbol{\theta}\|_1$  with  $\lambda=0.5\sqrt{\frac{\log p}{N}}$ , such that the regularized MLE over the whole dataset recovers the nonzeros of  $\boldsymbol{\theta}^*$  accurately. Figure 2 shows the performance of CEASE algorithms and ADMM, where "large n", "moderate n" and "small n" refer to (n,m)=(1000,5), (500,10) and (250,20), and "zero initialization" and "good initialization" refer to  $\boldsymbol{\theta}_0=\mathbf{0}$  and  $\bar{\boldsymbol{\theta}}$ , respectively. Again,  $\bar{\boldsymbol{\theta}}$  is the one-shot distributed estimator (Zhang et al., 2013). All the results are average values of 100 independent runs.

Again, the CEASE algorithms with proper regularization (Algorithms 3 and 4) work well in general; without regularization, the CEASE algorithm fails to converge when the local sample size n is small and the initialization in uninformative. For this nonsmooth problem,

the CEASE algorithm with averaging (Algorithm 4) does not seem to have advantage over the single version (Algorithm 3). The ADMM converges quickly to a region near the minimizer but then proceeds quite slowly, which appears to be a common phenomenon in many distributed optimization problems (Boyd et al., 2011).

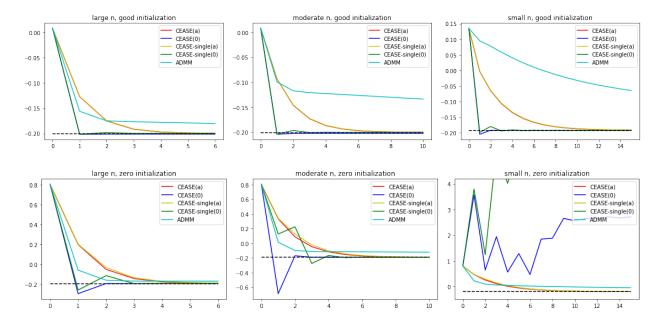


Figure 2: Nonsmooth minimization problems. The x-axis and y-axis are the number of iterations and  $\log \|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|_2$ . The dashed lines show the error of the minimizer of the overall regularized loss function. Top panel uses  $\bar{\boldsymbol{\theta}}$  as the initial value and bottom uses 0 as the initial value. CEASE(a) and CEASE(0) refer to Algorithm 4 with  $\alpha = 0.05p/n$  and 0; CEASE-single(a) and CEASE-single(0) refer to Algorithm 3 with  $\alpha = 0.05p/n$  and 0, respectively.

To summarize, the above simulations demonstrate several important properties of the CEASE Algorithms:

- In all scenarios, the CEASE Algorithms converge rapidly, usually within several steps; this is consistent with our theory.
- The CEASE Algorithms efficiently utilizes statistical structures and similarities among local losses, and benefit from the averaging step with smooth loss functions;
- The CEASE Algorithms are also able to handle the most general situations (e.g. small local sample size, uninformative initialization) with convergence guarantees.

#### 4.2 Real data

As a real data example, we choose the Spambase dataset from the UCI machine learning repository (Dua and Graff, 2017) as a testbed for comparison of algorithms. The goal is to train a classifier that distinguishes spam emails from normal ones, all of which are represented by 57-dimensional feature vectors based on their word frequencies and other characteristics. The total sample size is 4600. We build the testing set by randomly selecting 1000 samples from the entire dataset, and conducting logistic regression in a distributed manner on the rest of 3600 samples. We use the classification error on the testing set as a metric. Figure 3 shows the average performance of the CEASE algorithms, ADMM, GIANT and AGD based on 100 independent runs, where "large n", "moderate n" and "small n" refer to (n, m) = (720, 5), (360, 10) and (180, 20), respectively. All of the iterations are initialized with the one-shot average (Zhang et al., 2013). The experiments on this real data example also support our theoretical findings.

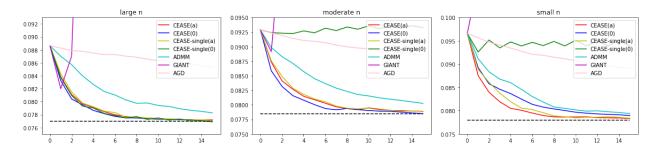


Figure 3: Spambase dataset. The x-axis and y-axis are the number of iterations and the testing error. The dashed lines show the error of the classifier based on all of the training samples. All of the iterations are initialized with the one-shot average  $\bar{\theta}$ . CEASE(a) and CEASE(0) refer to Algorithm 4 with  $\alpha = 0.15p/n$  and 0; CEASE-single(a) and CEASE-single(0) refer to Algorithm 3 with  $\alpha = 0.15p/n$  and 0, respectively. GIANT does not converge in all the three regimes, neither does CEASE-single(0) in the last two regimes.

# 5 Discussions

We have developed two CEASE distributed estimators (Algorithms 3 and 4) for statistical estimation, established their theoretical guarantees and carried out numerical experiments to illustrate their superior performance. Several new directions are worth exploring in the future. First, while we assumed exact computation in each step for simplicity, finer analysis

should allow for inexact updates and provide a guideline for practice. Second, we focus on the scenario where multiple node machines are all connected to a central processor, and all updates are done simultaneously. It would be interesting to generalize the algorithms to decentralized and asynchronous setting. Third, communication-efficient versions of confidence regions and hypothesis tests for sparse regression are of great importance in distributed statistical inference, and our point estimation strategies may serve as a starting point. It will be interesting to study how sparsity affects the contraction rates. Finally, it will be of great importance to explore non-convex statistical optimization problems such as mixture models and deep learning models. We believe that the idea of gradient-enhanced loss function still plays an important role to such an endeavor.

# A Proofs

## A.1 Proof of Theorem 2.1

Theorem 2.1 is a direct summary of the following two lemmas.

**Lemma A.1** (Contraction). Let Assumptions 2.1 and 2.2 hold, with  $\rho_0 > \delta \geq 0$ . Then  $\|\varphi_k(\boldsymbol{\theta}) - \widehat{\boldsymbol{\theta}}\|_2 \leq (\delta/\rho_0) \|\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}\|_2$ ,  $\forall \boldsymbol{\theta} \in B(\widehat{\boldsymbol{\theta}}, R), \forall k \in [m]$ .

*Proof.* Fix  $\boldsymbol{\theta} \in B(\widehat{\boldsymbol{\theta}}, R)$ . By the first order condition of  $\varphi_k(\boldsymbol{\theta})$ , we have that

$$\nabla f_k(\boldsymbol{\theta}) - \nabla f(\boldsymbol{\theta}) \in \partial \{ f_k[\varphi_k(\boldsymbol{\theta})] + g[\varphi_k(\boldsymbol{\theta})] \}. \tag{A.1}$$

Using the fixed point property  $\varphi_k(\widehat{\boldsymbol{\theta}}) = \widehat{\boldsymbol{\theta}}$ , we have  $\nabla f_k(\widehat{\boldsymbol{\theta}}) - \nabla f(\widehat{\boldsymbol{\theta}}) \in \partial [f_k(\widehat{\boldsymbol{\theta}}) + g(\widehat{\boldsymbol{\theta}})]$ . By the Taylor expansion and Assumption 2.2,

$$\begin{aligned} &\|[\nabla f_k(\boldsymbol{\theta}) - \nabla f(\boldsymbol{\theta})] - [\nabla f_k(\widehat{\boldsymbol{\theta}}) - \nabla f(\widehat{\boldsymbol{\theta}})]\|_2 \\ &= \left\| \int_0^1 \left( \nabla^2 f_k[(1-t)\widehat{\boldsymbol{\theta}} + t\boldsymbol{\theta}] - \nabla^2 f[(1-t)\widehat{\boldsymbol{\theta}} + t\boldsymbol{\theta}] \right) (\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}) dt \right\|_2 \\ &\leq \sup_{\boldsymbol{\zeta} \in B(\widehat{\boldsymbol{\theta}}, R)} \|\nabla^2 f_k(\boldsymbol{\zeta}) - \nabla^2 f(\boldsymbol{\zeta})\|_2 \cdot \|\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}\|_2 \\ &\leq \delta \|\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}\|_2 < \rho_0 R. \end{aligned}$$

From this, (A.1) and Lemma B.2, we obtain that  $\|\varphi_k(\boldsymbol{\theta}) - \widehat{\boldsymbol{\theta}}\|_2 \leq (\delta/\rho_0)\|\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}\|_2$ .

**Lemma A.2** (Averaging). Let Assumptions 2.1, 2.2 and 2.3 hold, with  $\rho_0 > \delta \geq 0$ . We have

$$\left\| \frac{1}{m} \sum_{k=1}^{m} \varphi_k(\boldsymbol{\theta}) - \widehat{\boldsymbol{\theta}} \right\|_2 \le \frac{\delta^2}{\rho_0 \rho} (1 + M \rho_0^{-1} \|\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}\|_2) \|\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}\|_2, \quad \forall \boldsymbol{\theta} \in B(\widehat{\boldsymbol{\theta}}, R).$$

Proof. Define  $L_k(\boldsymbol{\theta}) = f_k(\boldsymbol{\theta}) + g(\boldsymbol{\theta})$  and  $L(\boldsymbol{\theta}) = f(\boldsymbol{\theta}) + g(\boldsymbol{\theta})$  for  $\boldsymbol{\theta} \in \mathbb{R}^p$ . Then  $\widehat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\xi} \in \mathbb{R}^p} L(\boldsymbol{\xi})$  and  $\varphi_k(\boldsymbol{\theta}) = \operatorname{argmin}_{\boldsymbol{\xi} \in \mathbb{R}^p} \{L_k(\boldsymbol{\xi}) - \langle \nabla L_k(\boldsymbol{\theta}) - \nabla L(\boldsymbol{\theta}), \boldsymbol{\xi} \rangle \}$ . By the optimality conditions,

$$\nabla L_k[\varphi_k(\boldsymbol{\theta})] - \nabla L_k(\boldsymbol{\theta}) + \nabla L(\boldsymbol{\theta}) = \mathbf{0} = \nabla L(\widehat{\boldsymbol{\theta}}).$$

After subtracting  $\nabla L_k(\widehat{\boldsymbol{\theta}})$  from both sides and rearranging terms, we get

$$\nabla L_k[\varphi_k(\boldsymbol{\theta})] - \nabla L_k(\widehat{\boldsymbol{\theta}}) = [\nabla L_k(\boldsymbol{\theta}) - \nabla L_k(\widehat{\boldsymbol{\theta}})] - [\nabla L(\boldsymbol{\theta}) - \nabla L(\widehat{\boldsymbol{\theta}})].$$

Note that the average of the right hand side over  $k \in [m]$  is **0**.

Define 
$$\mathbf{H}_k = \int_0^1 \nabla^2 L_k[(1-t)\widehat{\boldsymbol{\theta}} + t\varphi_k(\boldsymbol{\theta})] dt$$
 for  $k \in [m]$  and  $\widehat{\mathbf{H}} = \nabla^2 L(\widehat{\boldsymbol{\theta}})$ . Then

$$\nabla L_k[\varphi_k(\boldsymbol{\theta})] - \nabla L_k(\widehat{\boldsymbol{\theta}}) = \mathbf{H}_k(\varphi_k(\boldsymbol{\theta}) - \widehat{\boldsymbol{\theta}}) = \widehat{\mathbf{H}}(\varphi_k(\boldsymbol{\theta}) - \widehat{\boldsymbol{\theta}}) + (\mathbf{H}_k - \widehat{\mathbf{H}})(\varphi_k(\boldsymbol{\theta}) - \widehat{\boldsymbol{\theta}}),$$

$$\mathbf{0} = \frac{1}{m} \sum_{k=1}^m \left( \nabla L_k[\varphi_k(\boldsymbol{\theta})] - \nabla L_k(\widehat{\boldsymbol{\theta}}) \right) = \widehat{\mathbf{H}}[\bar{\varphi}(\boldsymbol{\theta}) - \widehat{\boldsymbol{\theta}}] + \frac{1}{m} \sum_{k=1}^m (\mathbf{H}_k - \widehat{\mathbf{H}})(\varphi_k(\boldsymbol{\theta}) - \widehat{\boldsymbol{\theta}}),$$

where we let  $\bar{\varphi}(\boldsymbol{\theta}) = \frac{1}{m} \sum_{k=1}^{m} \varphi_k(\boldsymbol{\theta})$ . As a result,

$$\|\bar{\varphi}(\boldsymbol{\theta}) - \widehat{\boldsymbol{\theta}}\|_{2} = \left\| \frac{1}{m} \sum_{k=1}^{m} \widehat{\mathbf{H}}^{-1} (\mathbf{H}_{k} - \widehat{\mathbf{H}}) (\varphi_{k}(\boldsymbol{\theta}) - \widehat{\boldsymbol{\theta}}) \right\|_{2}$$

$$\leq \|\widehat{\mathbf{H}}^{-1}\|_{2} \max_{k \in [m]} \|\mathbf{H}_{k} - \widehat{\mathbf{H}}\|_{2} \cdot \max_{k \in [m]} \|\varphi_{k}(\boldsymbol{\theta}) - \widehat{\boldsymbol{\theta}}\|_{2}.$$

Lemma A.1 forces  $\max_{k \in [m]} \|\varphi_k(\boldsymbol{\theta}) - \widehat{\boldsymbol{\theta}}\|_2 \le (\delta/\rho_0) \|\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}\|_2$ , and Assumption 2.1 yields  $\widehat{\mathbf{H}} \succeq \rho \mathbf{I}$  and  $\|\widehat{\mathbf{H}}\|_2 \le 1/\rho$ . Furthermore, we use Assumptions 2.2 and 2.3 to get

$$\|\mathbf{H}_{k} - \widehat{\mathbf{H}}\|_{2} \leq \left\| \int_{0}^{1} \left( \nabla^{2} L_{k}[(1-t)\widehat{\boldsymbol{\theta}} + t\varphi_{k}(\boldsymbol{\theta})] - \nabla^{2} L[(1-t)\widehat{\boldsymbol{\theta}} + t\varphi_{k}(\boldsymbol{\theta})] \right) dt \right\|_{2}$$

$$+ \left\| \int_{0}^{1} \left( \nabla^{2} L[(1-t)\widehat{\boldsymbol{\theta}} + t\varphi_{k}(\boldsymbol{\theta})] - \nabla^{2} L(\widehat{\boldsymbol{\theta}}) \right) dt \right\|_{2}$$

$$\leq \delta + M \|\varphi_{k}(\boldsymbol{\theta}) - \widehat{\boldsymbol{\theta}}\|_{2} \leq \delta + M(\delta/\rho_{0}) \|\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}\|_{2}.$$

The proof is finished by combining all the estimates above.

## A.2 Proof of Theorem 2.2

Theorem 2.2 is a special of Theorem 3.2 by taking  $\alpha = 0$ . See section A.6 for proof of Theorem 3.2.

## A.3 Proof of Theorem 3.1

**Lemma A.3.** Let Assumptions 2.1 and 2.2 hold. Consider the iterates  $\{\boldsymbol{\theta}_t\}_{t=0}^T$  generated by Algorithm 4. Define

$$\gamma_t = \begin{cases} \frac{\delta}{\rho_0 + \alpha} \cdot \min\{1, \frac{\delta}{\rho + \alpha}(1 + \frac{M}{\rho_0 + \alpha} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2)\} &, \text{ if Assumption 2.3 holds} \\ \frac{\delta}{\rho_0 + \alpha} &, \text{ otherwise} \end{cases}, \quad 0 \le t \le T - 1.$$

If  $0 < \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2 < R/2$ ,  $\delta < \rho_0 + \alpha$  and  $\gamma_t^2 < \rho/(\rho + 2\alpha)$ , then

$$\frac{\|\boldsymbol{\theta}_{t+1} - \widehat{\boldsymbol{\theta}}\|_2}{\|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2} \le \frac{\gamma_t \sqrt{\rho^2 + 2\alpha\rho} + \alpha}{\rho + \alpha} < 1.$$

Theorem 3.1 directly follows from Lemma A.3 and induction. Below we only prove Lemma A.3 with Assumption 2.3. The other part in Lemma A.3 without Assumption 2.3 can be derived by slightly modifying this proof.

Proof of Lemma A.3 with Assumption 2.3. Let  $\theta_t^+ = \text{prox}_{\alpha^{-1}(f+g)}(\theta_t)$ . By the triangle inequality,

$$\|\boldsymbol{\theta}_{t+1} - \widehat{\boldsymbol{\theta}}\|_{2} \le \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t}^{+}\|_{2} + \|\boldsymbol{\theta}_{t}^{+} - \widehat{\boldsymbol{\theta}}\|_{2}.$$
 (A.2)

We first invoke Theorem 2.1 to bound the first term  $\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t^+\|_2$  in (A.2). Define  $\tilde{g}(\boldsymbol{\theta}) = g(\boldsymbol{\theta}) + (\alpha/2)\|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|_2^2$  for  $\boldsymbol{\theta} \in \mathbb{R}^p$ . Then  $\boldsymbol{\theta}_{t+1}$  is the first iterate of Algorithm 2 initialized at  $\boldsymbol{\theta}_t$  for computing  $\boldsymbol{\theta}_t^+ = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^p} \{ \frac{1}{m} \sum_{k=1}^m f_k(\boldsymbol{\theta}) + \tilde{g}(\boldsymbol{\theta}) \}$ .

From  $\widehat{\boldsymbol{\theta}} = \operatorname{prox}_{\alpha^{-1}(f+g)}(\widehat{\boldsymbol{\theta}})$  and Lemma B.3 we obtain that  $\|\boldsymbol{\theta}_t^+ - \widehat{\boldsymbol{\theta}}\|_2 \leq \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2$ . Then the condition  $\|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2 < R/2$  leads to  $B(\boldsymbol{\theta}_t^+, R/2) \subseteq B(\widehat{\boldsymbol{\theta}}, R)$ . By Assumptions 2.1 and 2.2,

- in  $B(\boldsymbol{\theta}_t^+, R/2)$ ,  $\{f_k + \tilde{g}\}_{k=1}^m$  are  $(\rho_0 + \alpha)$ -strongly convex and  $(f + \tilde{g})$  is  $(\rho + \alpha)$ -strongly convex;
- $\|\nabla^2 f_k(\boldsymbol{\theta}) \nabla^2 f(\boldsymbol{\theta})\|_2 \le \delta$  holds for all  $k \in [m]$  and  $\boldsymbol{\theta} \in B(\boldsymbol{\theta}_t^+, R/2)$ .

Furthermore, Lemma B.3 also yields

$$\|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}_{t}^{+}\|_{2}^{2} \leq \|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2}^{2} - \|\boldsymbol{\theta}_{t}^{+} - \widehat{\boldsymbol{\theta}}\|_{2}^{2} = \|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2}^{2} \left(1 - \|\boldsymbol{\theta}_{t}^{+} - \widehat{\boldsymbol{\theta}}\|_{2}^{2} / \|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2}^{2}\right). \tag{A.3}$$

Then  $\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_t^+\|_2 \le \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2 < R/2$  and  $\tilde{\boldsymbol{\theta}}_0 \in B(\boldsymbol{\theta}_t^+, R/2)$ . Based on these conditions and  $\alpha + \rho_0 > \delta$ , we use Theorem 2.1 to get

$$\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t^+\|_2 \le \gamma_t \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_t^+\|_2.$$

From here, (A.2) and (A.3) we obtain that

$$\begin{aligned} \|\boldsymbol{\theta}_{t+1} - \widehat{\boldsymbol{\theta}}\|_{2} &\leq \gamma_{t} \|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}_{t}^{+}\|_{2} + \|\boldsymbol{\theta}_{t}^{+} - \widehat{\boldsymbol{\theta}}\|_{2} \\ &\leq \gamma_{t} \|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2} \left(1 - \|\boldsymbol{\theta}_{t}^{+} - \widehat{\boldsymbol{\theta}}\|_{2}^{2} / \|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2}^{2}\right)^{1/2} + \|\boldsymbol{\theta}_{t}^{+} - \widehat{\boldsymbol{\theta}}\|_{2} \\ &= \|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2} \cdot h(\|\boldsymbol{\theta}_{t}^{+} - \widehat{\boldsymbol{\theta}}\|_{2} / \|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2}), \end{aligned}$$

where  $h(x) = \gamma_t \sqrt{1 - x^2} + x$ ,  $\forall x \in [0, 1]$ . From  $h'(x) = 1 - \gamma_t x / \sqrt{1 - x^2}$  we see that  $h' \ge 0$  on  $[0, 1/\sqrt{1 + \gamma_t^2}]$ .

On the one hand, Lemma B.3 asserts that  $\|\boldsymbol{\theta}^+ - \widehat{\boldsymbol{\theta}}\|_2 / \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2 \le \alpha / (\rho + \alpha)$ . On the other hand, the assumption  $\gamma_t^2 < \rho / (\rho + 2\alpha)$  forces

$$\frac{1}{\sqrt{1+\gamma_t^2}} > \frac{1}{\sqrt{1+\rho/(\rho+2\alpha)}} = \frac{\sqrt{\rho/2+\alpha}}{\sqrt{\rho+\alpha}} \ge \frac{\rho/2+\alpha}{\rho+\alpha} \ge \alpha/(\rho+\alpha).$$

The proof is completed by computation:

$$\frac{\|\boldsymbol{\theta}_{t+1} - \widehat{\boldsymbol{\theta}}\|_{2}}{\|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2}} \leq h \left( \frac{\|\boldsymbol{\theta}_{t}^{+} - \widehat{\boldsymbol{\theta}}\|_{2}}{\|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2}} \right) \leq h \left( \frac{\alpha}{\rho + \alpha} \right) = \gamma_{t} \left[ 1 - \left( \frac{\alpha}{\rho + \alpha} \right)^{2} \right]^{1/2} + \frac{\alpha}{\rho + \alpha}$$

$$= \frac{\gamma_{t} \sqrt{\rho^{2} + 2\rho\alpha} + \alpha}{\rho + \alpha} < \frac{\sqrt{[\rho/(\rho + 2\alpha)] \cdot (\rho^{2} + 2\rho\alpha)} + \alpha}{\rho + \alpha} = 1,$$

where we used the assumption  $\gamma_t^2 < \rho/(\rho + 2\alpha)$  again.

# A.4 Proof of Corollary 3.1

We claim that  $(\frac{\delta}{\rho_0 + \alpha})^2 \leq \frac{7}{9} \cdot \frac{\rho}{\rho + 2\alpha}$ . Given this, Corollary 3.1 follows from Theorem 3.1 and

$$\frac{\frac{\delta}{\rho_0 + \alpha} \sqrt{\rho^2 + 2\alpha\rho} + \alpha}{\rho + \alpha} \le \frac{\sqrt{\frac{7}{9} \cdot \frac{\rho}{\rho + 2\alpha} \cdot \rho(\rho + 2\alpha)} + \alpha}{\rho + \alpha} = 1 - \frac{(1 - \sqrt{7/9})\rho}{\rho + \alpha} \le 1 - \frac{\rho/10}{\rho + \alpha}.$$

The claim trivially holds if  $\delta = 0$ . When  $\delta > 0$ , let us first assume  $0 < \delta \le \rho$  and define

 $b = \alpha \rho / \delta^2$ . Then  $\alpha = b \delta^2 / \rho$ ,  $b \ge 4$  and  $\rho_0 \ge \max \{ \rho - \delta, 0 \}$  force

$$\rho_0 + \alpha \ge \rho - \delta + b\delta^2/\rho = (\rho/\delta + b\delta/\rho - 1)\delta \ge (2\sqrt{(\rho/\delta) \cdot (b\delta/\rho)} - 1)\delta = (2\sqrt{b} - 1)\delta > \delta.$$

and  $[\delta/(\rho_0 + \alpha)]^2 \leq [\delta/(\rho_0 + \alpha)]^2 \leq 1/h_1(\delta/\rho)$ , where  $h_1(x) = (bx + x^{-1} - 1)^2$ . On the other hand,  $\rho/(\rho + 2\alpha) = 1/(1 + 2\alpha/\rho) = 1/h_2(\delta/\rho)$ , where  $h_2(x) = 1 + 2bx^2$ .

We are going to show  $h_2(x) \leq 7h_1(x)/9$ ,  $\forall x \in (0,1]$ , which leads to the desired result under  $0 < \delta \leq \rho$ . If  $0 < x \leq \sqrt{3}/2$ , then  $h_1(x) \geq (2\sqrt{b} - 1)^2 \geq (2\sqrt{b} - \sqrt{b}/2)^2 \geq 9b/4$  and

$$h_2(x) \le 1 + 2b \cdot (3/4) \le (b/4) + (6b/4) = 7b/4 \le 7h_1(x)/9.$$

If  $\sqrt{3}/2 < x \le 1$ , then  $h_1(x) \ge b^2 x^2 \ge 3b^2/4$ ,  $h_2(x) \le 1 + 2b \le (b/4) + 2b = 9b/4$ , and  $h_2(x)/h_1(x) = 3/b \le 3/4 \le 7/9$ .

Suppose now that  $\delta > \rho$ , and define  $b = \alpha \rho / \delta^2$ . Then

$$\left(\frac{\delta}{\rho_0 + \alpha}\right)^2 \le \left(\frac{\delta}{\alpha}\right)^2 = \left(\frac{1}{b\delta/\rho}\right)^2 = \frac{1}{b^2(\delta/\rho)^2},$$
$$\frac{\rho}{\rho + 2\alpha} = \frac{1}{1 + 2\alpha/\rho} = \frac{1}{1 + 2b(\delta/\rho)^2}.$$

From  $b \geq 4$  and  $\delta/\rho > 1$  we get  $(\frac{\delta}{\rho_0 + \alpha})^2 \leq \frac{7}{9} \cdot \frac{\rho}{\rho + 2\alpha}$  from

$$b^{2}(\delta/\rho)^{2} - \frac{9}{7}[1 + 2b(\delta/\rho)^{2}] = (\delta/\rho)^{2}b(b - 18/7) - 9/7$$
  
 
$$\geq 1 \cdot 4 \cdot (4^{1} - 18/7) - 9/7 = 31/7 > 0.$$

# A.5 Proof of Corollary 3.2

Throughout the proof we assume that  $\delta/\rho$  is sufficiently small. The regularity conditions in Theorem 3.1 are easily verified as  $\boldsymbol{\theta}_0 \in B(\widehat{\boldsymbol{\theta}}, R/2)$  and  $[\delta/(\rho_0 + \alpha)]^2 < \rho/(\rho + 2\alpha)$ . Here we used the fact  $\rho_0 \ge \rho - \delta$ .

From  $\rho_0 \ge \rho - \delta$  we get  $\rho_0 + \alpha \ge \rho_0 \ge \rho/2$  and  $\delta/(\rho_0 + \alpha) \le 2\delta/\rho$ . Also,  $\sqrt{\rho^2 + 2\alpha\rho} \le \rho\sqrt{1 + 2C(\delta/\rho)^2} \lesssim \rho$ . We control the contraction factor in (3.2):

$$\frac{\frac{\delta}{\rho_0 + \alpha} \sqrt{\rho^2 + 2\alpha\rho} + \alpha}{\rho + \alpha} \lesssim \frac{(2\delta/\rho)\rho + C\delta^2/\rho}{\rho} = \frac{2\delta}{\rho} + \frac{C}{4} \left(\frac{2\delta}{\rho}\right)^2 \lesssim \frac{\delta}{\rho}.$$

Recall that  $\gamma_t = \frac{\delta}{\rho_0 + \alpha} \cdot \min\{1, \frac{\delta}{\rho + \alpha}(1 + \frac{M}{\rho_0 + \alpha} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2)\}$  in Theorem 3.1. When  $\delta/\rho$  is small and  $\|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2 \le \rho/M$ , we have  $\rho_0 + \alpha \ge \rho/2$ ,  $\frac{M}{\rho_0 + \alpha} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2 \le 2$ , and  $\gamma_t \le (2\delta/\rho)^2$ . This help bound the contraction factor in (3.2):

$$\frac{\gamma_t \cdot \sqrt{\rho^2 + 2\alpha\rho} + \alpha}{\rho + \alpha} \lesssim \frac{(2\delta/\rho)^2 \rho + C\delta^2/\rho}{\rho} \lesssim \left(\frac{\delta}{\rho}\right)^2.$$

## A.6 Proof of Theorem 3.2

The proof is implied by combining proof of Corollary 3.2 with the results of the following two lemmas, the first of which is a direct counterpart of Theorem 3.1 in the stochastic setting, given an additional condition on similarity between local Hessians. The second lemma below specifies the order of Hessian difference in the generalized linear model, hence providing a contraction rate and guiding the choice of  $\alpha$ .

Lemma A.4. Let Assumption 2.5 hold. Denote

$$\widehat{\delta} := 2 \sup_{k \in [m]} \sup_{\boldsymbol{\theta} \in B(\widehat{\boldsymbol{\theta}}, R)} \|\nabla^2 f_k(\boldsymbol{\theta}) - \nabla^2 F(\boldsymbol{\theta})\|_2.$$

Consider the iterates  $\{\boldsymbol{\theta}_t\}_{t=0}^T$  generated by Algorithm 3 or Algorithm 4. Suppose that  $\boldsymbol{\theta}_0 \in B(\widehat{\boldsymbol{\theta}}, R/2)$  and  $[\widehat{\delta}/(\rho_0 + \alpha)]^2 < \rho/(\rho + 2\alpha)$ .

• For both Algorithms 3 and 4, we have

$$\|\boldsymbol{\theta}_{t+1} - \widehat{\boldsymbol{\theta}}\|_{2} \leq \|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2} \cdot \frac{\widehat{\delta}}{\frac{\rho_{0} + \alpha}{\rho_{0} + \alpha}} \cdot \sqrt{\rho^{2} + 2\alpha\rho} + \alpha}{\rho + \alpha}, \qquad 0 \leq t \leq T - 1; \tag{A.4}$$

• If in addition, Assumption 2.6 also holds, then for Algorithm 4 we have

$$\|\boldsymbol{\theta}_{t+1} - \widehat{\boldsymbol{\theta}}\|_{2} \le \|\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\|_{2} \cdot \frac{\gamma_{t}\sqrt{\rho^{2} + 2\alpha\rho} + \alpha}{\rho + \alpha}, \qquad 0 \le t \le T - 1,$$
 (A.5)

where we define  $\gamma_t = \frac{\hat{\delta}}{\rho_0 + \alpha} \cdot \min\{1, \frac{\hat{\delta}}{\rho + \alpha}(1 + \frac{M}{\rho_0 + \alpha} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}\|_2)\};$ 

• Both multiplicative factors in (A.4) and (A.5) are strictly less than 1.

**Proof of Lemma A.4.** We first assume that  $\rho_0 > \widehat{\delta}$  and analyze the vanilla DANE algorithm under the new assumptions. Let Assumption 2.5 hold, and  $\{\boldsymbol{\theta}_t\}_{t=0}^{\infty}$  be the iterates with  $\boldsymbol{\theta}_0 \in B(\widehat{\boldsymbol{\theta}}, R)$ . For any  $\boldsymbol{\theta} \in B(\widehat{\boldsymbol{\theta}}, R)$ ,  $\|\nabla^2 f(\boldsymbol{\theta}) - \nabla^2 F(\boldsymbol{\theta})\|_2 \leq \widehat{\delta}/2$  and thus

 $\|\nabla^2 f_k(\boldsymbol{\theta}) - \nabla^2 f(\boldsymbol{\theta})\|_2 \leq \hat{\delta}$  for  $k \in [m]$ . Hence it implies Assumption 2.2 with  $\delta = \hat{\delta}$ , and Lemma A.1 continues to hold. We can also get the result in Lemma A.2 under Assumption 2.6, by replacing  $\hat{\mathbf{H}} = \nabla^2 f(\hat{\boldsymbol{\theta}})$  in the proof of Lemma A.2 by  $\nabla^2 F(\hat{\boldsymbol{\theta}}) + \nabla^2 g(\boldsymbol{\theta})$ . Then we drop the assumption  $\rho_0 > \hat{\delta}$  can reproduce the results in Theorem 3.1 under the new setting, by following its original proof.

**Lemma A.5.** Under Assumption 2.4, for an arbitrarily small positive constant c, there exist universal constants  $C_1, C_2$  and  $C_3$  depending only on c such that as long as  $n \ge cp$ , with probability at least  $1 - 2e^{-C_2n} - Ne^{-C_3p}$ ,

$$\sup_{k \in [m]} \sup_{\boldsymbol{\theta} \in B(\widehat{\boldsymbol{\theta}}, R)} \|\nabla^2 f_k(\boldsymbol{\theta}) - \nabla^2 F(\boldsymbol{\theta})\|_2 \le C_1 \|\boldsymbol{\Sigma}\|_2 \sqrt{\frac{p \max\{1, \log(Np^{1/2} \|\boldsymbol{\Sigma}\|_2 R)\}}{n}}.$$

**Proof of Lemma A.5.** Let  $\tilde{\mathbf{x}}_i = (\mathbf{\Sigma}^*)^{-1/2} \mathbf{x}_i$ ,  $\tilde{\boldsymbol{\theta}} = (\mathbf{\Sigma}^*)^{1/2} \boldsymbol{\theta}$ , and define a new loss function  $\tilde{l}(\tilde{\boldsymbol{\theta}}, \tilde{\mathbf{x}}_i) = b(\tilde{\mathbf{x}}_i^{\top} \tilde{\boldsymbol{\theta}}) - y_i(\tilde{\mathbf{x}}_i^{\top} \tilde{\boldsymbol{\theta}})$ . Let  $\hat{R}_k(\tilde{\boldsymbol{\theta}}) = \frac{1}{n} \sum_{i \in \mathcal{I}_k} \tilde{l}(\tilde{\boldsymbol{\theta}}, \tilde{\mathbf{x}}_i)$  for  $k \in [m]$ . Then we have that  $\nabla^2 \hat{R}_k(\tilde{\boldsymbol{\theta}}) = \frac{1}{n} \sum_{i \in \mathcal{I}_k} b''(\tilde{\boldsymbol{\theta}}^{\top} \tilde{\mathbf{x}}_i) \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^{\top}$  and  $\nabla^2 f_k(\boldsymbol{\theta}) = (\mathbf{\Sigma}^*)^{1/2} \nabla^2 \hat{R}_k(\tilde{\boldsymbol{\theta}}) (\mathbf{\Sigma}^*)^{1/2}$ . Similarly we have  $\nabla^2 F(\boldsymbol{\theta}) = (\mathbf{\Sigma}^*)^{1/2} \mathbb{E} \nabla^2 \hat{R}_k(\tilde{\boldsymbol{\theta}}) (\mathbf{\Sigma}^*)^{1/2}$ .

Therefore

$$\max_{k \in [m]} \max_{\boldsymbol{\theta} \in B(\widehat{\boldsymbol{\theta}}, R)} \|\nabla^2 f_k(\boldsymbol{\theta}) - \nabla^2 F(\boldsymbol{\theta})\|_2 \le \|\boldsymbol{\Sigma}^*\|_2 \max_{k \in [m]} \max_{\tilde{\boldsymbol{\theta}} \in B((\boldsymbol{\Sigma}^*)^{1/2}\widehat{\boldsymbol{\theta}}, \tilde{R})} \|\nabla^2 \widehat{R}_k(\tilde{\boldsymbol{\theta}}) - \nabla^2 \mathbb{E} \widehat{R}_k(\tilde{\boldsymbol{\theta}})\|_2$$
(A.6)

and we only need to control the quantity on the right hand side. Here  $\tilde{R} = \|\mathbf{\Sigma}^*\|_2^{1/2} R$ .

Define  $\Delta_0 = \max_{k \in [m]} \max_{\tilde{\boldsymbol{\theta}} \in B((\boldsymbol{\Sigma}^*)^{1/2} \hat{\boldsymbol{\theta}}, \tilde{R})} \|\nabla^2 \hat{R}_k(\tilde{\boldsymbol{\theta}}) - \nabla^2 \mathbb{E} \hat{R}_k(\tilde{\boldsymbol{\theta}})\|_2$  and

$$\phi_k(\tilde{\boldsymbol{\theta}}) = \|\frac{1}{n} \sum_{i \in \mathcal{I}_k} b''(\tilde{\mathbf{x}}_i^{\top} \tilde{\boldsymbol{\theta}}) \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^{\top} - \mathbb{E}b''(\tilde{\mathbf{X}}^{\top} \tilde{\boldsymbol{\theta}}) \tilde{\mathbf{X}} \tilde{\mathbf{X}}^{\top} \|_2.$$

Here  $\tilde{\mathbf{X}}$  shares the distribution with  $\tilde{\mathbf{x}}_i$ . Firstly we bound  $\phi_k(\tilde{\boldsymbol{\theta}})$  for any fixed  $\tilde{\boldsymbol{\theta}} \in B(\tilde{\boldsymbol{\theta}}^*, 2\tilde{R})$ , where  $\tilde{\boldsymbol{\theta}}^* := (\boldsymbol{\Sigma}^*)^{1/2} \hat{\boldsymbol{\theta}}$ . For any  $k \in [m]$ , under Assumption 2.4, there exist constants  $c_1, c_2$  such that for any  $\epsilon \geq 0$ 

$$\mathbb{P}(\phi_k(\tilde{\boldsymbol{\theta}}) \ge \epsilon) \le 2e^{c_1 p - c_2 \min\{\epsilon, \epsilon^2\} n}. \tag{A.7}$$

To see this, notice that  $\phi_k(\tilde{\boldsymbol{\theta}}) = \sup_{\mathbf{u} \in \mathbb{S}^{p-1}} g_k(\mathbf{u})$ , where  $g_k(\mathbf{u}) = \mathbf{u}^\top \{\frac{1}{n} \sum_{i \in \mathcal{I}_k} b''(\tilde{\mathbf{x}}_i^\top \tilde{\boldsymbol{\theta}}) \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top - \mathbb{E}b''(\tilde{\mathbf{X}}^\top \tilde{\boldsymbol{\theta}}) \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top \} \mathbf{u}$ . Let  $\mathcal{N}$  be a  $\frac{1}{4}$ -covering of  $\mathbb{S}^{p-1}$ , and  $|\mathcal{N}| \leq 9^p$ . Denote  $\hat{\mathbf{u}} = \arg \max_{\mathbf{u}} g_k(\mathbf{u})$ .

Find  $\tilde{\mathbf{u}} \in \mathcal{N}$  such that  $\|\hat{\mathbf{u}} - \tilde{\mathbf{u}}\|_2 \leq \frac{1}{4}$ . Then

$$|g_k(\tilde{\mathbf{u}}) - g_k(\hat{\mathbf{u}})| = |(\tilde{\mathbf{u}} + \hat{\mathbf{u}})^\top \left\{ \frac{1}{n} \sum_{i \in \mathcal{I}_k} b''(\tilde{\mathbf{x}}_i^\top \tilde{\boldsymbol{\theta}}) \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top - \mathbb{E}b''(\tilde{\mathbf{X}}^\top \tilde{\boldsymbol{\theta}}) \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top \right\} (\tilde{\mathbf{u}} - \hat{\mathbf{u}})| \le \frac{1}{2} g_k(\hat{\mathbf{u}}),$$

and thus

$$\sup_{\mathbf{u} \in \mathbb{S}^{p-1}} g_k(\mathbf{u}) \le 2 \sup_{\mathbf{u} \in \mathcal{N}} g_k(\mathbf{u}).$$

On the other hand from Bernstein's inequality we see that there exists a constant  $c_2$  such that for any  $\mathbf{u} \in \mathcal{N}$ ,  $\epsilon \geq 0$ ,  $\mathbb{P}(g_k(\mathbf{u}) \geq \frac{\epsilon}{2}) \leq 2e^{-c_2 \min\{\epsilon, \epsilon^2\}n}$ . Therefore

$$\mathbb{P}(\phi_k(\tilde{\boldsymbol{\theta}}) \geq \epsilon) = \mathbb{P}(\sup_{\mathbf{u} \in \mathbb{S}^{p-1}} g_k(\mathbf{u}) \geq \epsilon) \leq \mathbb{P}(\sup_{\mathbf{u} \in \mathcal{N}} g_k(\mathbf{u}) \geq \frac{\epsilon}{2}) \leq |\mathcal{N}| \cdot 2e^{-c_2 \min\{\epsilon, \epsilon^2\}n} \leq 2e^{c_1 p - c_2 \min\{\epsilon, \epsilon^2\}n}$$

where  $c_1 = \log 9$ .

Now for  $t \geq 1$ , define the event  $E_t \triangleq \left\{ \max_{i=1}^N \|\tilde{\mathbf{x}}_i\|_2^3 < (8t)^{3/2} \mathbb{E} \|\tilde{\mathbf{X}}\|_2^3 \right\}$ . Then by Theorem 2.1 in Hsu et al. (2012),  $\mathbb{P}(E_t^c) = \mathbb{P}(\max_{i=1}^N \|\tilde{\mathbf{x}}_i\|_2^2 \geq 8t \mathbb{E} \|\tilde{\mathbf{X}}\|_2^3)^{2/3}) \leq \mathbb{P}(\max_{i=1}^N \|\tilde{\mathbf{x}}_i\|_2^2 \geq 8t \mathbb{E} \|\tilde{\mathbf{X}}\|_2^2) \leq N\mathbb{P}(\|\tilde{\mathbf{X}}\|_2^2 \geq 8t \mathbb{E} \|\tilde{\mathbf{X}}\|_2^2) \leq Ne^{-tp}$ . Under the event  $E_t$ , for  $\tilde{\boldsymbol{\theta}}_1, \tilde{\boldsymbol{\theta}}_2 \in B(\tilde{\boldsymbol{\theta}}^*, 2\tilde{R})$ , we have

$$|\phi_{k}(\tilde{\boldsymbol{\theta}}_{1}) - \phi_{k}(\tilde{\boldsymbol{\theta}}_{2})| \leq \|\frac{1}{n} \sum_{i \in \mathcal{I}_{k}} [b''(\tilde{\mathbf{x}}_{i}^{\top} \tilde{\boldsymbol{\theta}}_{1}) - b''(\tilde{\mathbf{x}}_{i}^{\top} \tilde{\boldsymbol{\theta}}_{2})] \tilde{\mathbf{x}}_{i} \tilde{\mathbf{x}}_{i}^{\top} \|_{2} + \|\mathbb{E}b''(\tilde{\mathbf{X}}^{\top} \tilde{\boldsymbol{\theta}}_{1}) \tilde{\mathbf{X}} \tilde{\mathbf{X}}^{\top} - \mathbb{E}b''(\tilde{\mathbf{X}}^{T} \tilde{\boldsymbol{\theta}}_{2}) \tilde{\mathbf{X}} \tilde{\mathbf{X}}^{\top} \|_{2}$$

$$\leq B_{3} \|\tilde{\boldsymbol{\theta}}_{1} - \tilde{\boldsymbol{\theta}}_{2} \|_{2} \cdot (\mathbb{E} \|\tilde{\mathbf{X}}\|_{2}^{3} + \frac{1}{n} \sum_{i=1}^{n} \|\tilde{\mathbf{x}}_{ki}\|_{2}^{3})$$

$$\leq (9t)^{3/2} B_{3} \cdot \mathbb{E} \|\mathbf{U}\|_{2}^{3} \cdot \|\tilde{\boldsymbol{\theta}}_{1} - \tilde{\boldsymbol{\theta}}_{2}\|_{2}$$

$$\leq c_{3}(pt)^{3/2} \|\tilde{\boldsymbol{\theta}}_{1} - \tilde{\boldsymbol{\theta}}_{2}\|_{2}$$

for some constant  $c_3$  depending only on  $B_3$ .

Now let  $\mathcal{N}_{\delta}$  be a  $\delta$ -covering of  $B(\tilde{\boldsymbol{\theta}}^*, 2\tilde{R})$ , where  $\delta = \frac{\epsilon}{c_3(tp)^{3/2}}$ . We can also assume that  $|\mathcal{N}_{\delta}| \leq (\frac{6\tilde{R}}{\delta})^p$ . Therefore for any  $k \in [m]$ ,

$$\mathbb{P}\left(E_t \cap \left\{\sup_{\tilde{\boldsymbol{\theta}} \in B(\tilde{\boldsymbol{\theta}}^*, 2\tilde{R})} \phi_k(\tilde{\boldsymbol{\theta}}) \geq 2\epsilon\right\}\right) \leq \mathbb{P}\left(E_t \cap \left\{\sup_{\tilde{\boldsymbol{\theta}} \in \mathcal{N}_{\delta}} \phi_k(\tilde{\boldsymbol{\theta}}) \geq \epsilon\right\}\right) \\
\leq |\mathcal{N}_{\delta}| \cdot 2e^{c_1 p - c_2 \min\{\epsilon, \epsilon^2\}n} = 2e^{c_4 p + c_3 p \log \frac{(tp)^{3/2}\tilde{R}}{\epsilon} - c_2 \min\{\epsilon, \epsilon^2\}n}.$$

Thus

$$\mathbb{P}(\boldsymbol{\Delta}_0 \ge 2\epsilon) \le \mathbb{P}\left(\bigcup_{k \in [m]} \left\{ \sup_{\tilde{\boldsymbol{\theta}} \in B(\tilde{\boldsymbol{\theta}}^*, \tilde{R})} \phi_k(\tilde{\boldsymbol{\theta}}) \ge 2\epsilon \right\} \right)$$

$$\le \mathbb{P}(E_t^c) + \sum_{k=1}^m \mathbb{P}\left(E_t \cap \left\{ \sup_{\tilde{\boldsymbol{\theta}} \in B(\tilde{\boldsymbol{\theta}}^*, \tilde{R})} \phi_k(\tilde{\boldsymbol{\theta}}) \ge 2\epsilon \right\} \right)$$

$$\le Ne^{-tp} + 2me^{c_4p + c_3p \log \frac{(tp)^{3/2}\tilde{R}}{\epsilon} - c_2 \min\{\epsilon, \epsilon^2\}n}$$

It is easily seen that the last expression is no more than  $Ne^{-tp} + 2e^{-C_2p}$  when

$$\begin{cases} \min\{\epsilon, \epsilon^2\} \ge C_1' \frac{\max\{p, p \log(t^{3/2}R), \log m\}}{n}, \\ \epsilon \ge 1 \text{ or } \frac{\epsilon^2}{\log \frac{1}{\epsilon}} \ge C_2' \cdot \frac{p}{n}, \end{cases}$$

which is satisfied if t is chosen as a suitable constant depending on c, and that

$$\epsilon = C\sqrt{\frac{\log m + p \max\{1, \log(np^{1/2}\tilde{R})\}}{n}}.$$
(A.8)

Here  $C'_i$  and C is a constant depending only on c.

Finally, note that

$$\Sigma^* = \operatorname{cov}(\mathbf{x}_i) = \begin{pmatrix} 1 \\ \Sigma \end{pmatrix}$$

and that  $\|\Sigma\|_2 \ge A_1$  for a universal  $A_1 > 0$ , we have  $\|\Sigma^*\|_2 \le \max\{1, 1/A_1\}\|\Sigma\|_2$ . Thus combining (A.8) and (A.6) completes the proof.

#### A.7 Proof of Theorem 3.3

We first present three lemmas, based on which we build the proof of the main theorem.

**Lemma A.6.** Suppose that  $\widehat{\Sigma}$  is positive-definite, i.e.  $\lambda_{\min}(\widehat{\Sigma}) > 0$ . Define  $\varepsilon_t = \widehat{\Sigma}^{1/2}(\boldsymbol{\theta}_t - \widehat{\Sigma}^{-1}\widehat{\mathbf{w}})$  for  $t \geq 0$  and

$$\mathbf{\Delta}_k = (\widehat{\mathbf{\Sigma}} + \alpha \mathbf{I})^{-1/2} (\widehat{\mathbf{\Sigma}}_k - \widehat{\mathbf{\Sigma}}) (\widehat{\mathbf{\Sigma}} + \alpha \mathbf{I})^{-1/2}, \quad \forall k \in [m].$$

If  $\alpha \geq 0$  is appropriately chosen such that  $\Delta = \max_{k \in [m]} \|\Delta_k\|_2 \leq 1/2$ . Then

$$\|\boldsymbol{\varepsilon}_{t+1}\|_{2} \leq \frac{2\Delta^{2} + \alpha/\lambda_{\min}(\widehat{\boldsymbol{\Sigma}})}{1 + \alpha/\lambda_{\min}(\widehat{\boldsymbol{\Sigma}})} \|\boldsymbol{\varepsilon}_{t}\|_{2}, \quad \forall t \geq 0,$$

which guarantees linear convergence of  $\{\varepsilon_t\}_{t=0}^{\infty}$ .

Proof of Lemma A.6. Define  $\boldsymbol{\varepsilon}_{t,k} = \widehat{\boldsymbol{\Sigma}}^{1/2} (\boldsymbol{\theta}_{t,k} - \widehat{\boldsymbol{\Sigma}}^{-1} \widehat{\mathbf{w}})$ . Then

$$\varepsilon_{t+1,k} = \widehat{\Sigma}^{1/2} (\boldsymbol{\theta}_{t+1,k} - \widehat{\Sigma}^{-1} \widehat{\mathbf{w}}) 
= \widehat{\Sigma}^{1/2} [\mathbf{I} - (\widehat{\Sigma}_k + \alpha \mathbf{I})^{-1} \widehat{\Sigma}] \boldsymbol{\theta}_t + \widehat{\Sigma}^{1/2} (\widehat{\Sigma}_k + \alpha \mathbf{I})^{-1} \widehat{\mathbf{w}} - \widehat{\Sigma}^{-1/2} \widehat{\mathbf{w}} 
= [\mathbf{I} - \widehat{\Sigma}^{1/2} (\widehat{\Sigma}_k + \alpha \mathbf{I})^{-1} \widehat{\Sigma}^{1/2}] \varepsilon_t.$$

Define  $\widetilde{\Sigma}_k^{(1)} = \widehat{\Sigma}^{-1/2} \widehat{\Sigma}_k \widehat{\Sigma}^{-1/2}$  for  $k \in [m]$ . The fact

$$\widehat{\boldsymbol{\Sigma}}_k + \alpha \mathbf{I} = \widehat{\boldsymbol{\Sigma}}^{1/2} (\widehat{\boldsymbol{\Sigma}}^{-1/2} \widehat{\boldsymbol{\Sigma}}_k \widehat{\boldsymbol{\Sigma}}^{-1/2} + \alpha \widehat{\boldsymbol{\Sigma}}^{-1}) \widehat{\boldsymbol{\Sigma}}^{1/2} = \widehat{\boldsymbol{\Sigma}}^{1/2} (\widehat{\boldsymbol{\Sigma}}_k^{(1)} + \alpha \widehat{\boldsymbol{\Sigma}}^{-1}) \widehat{\boldsymbol{\Sigma}}^{1/2}$$

gives 
$$\widehat{\boldsymbol{\Sigma}}^{1/2}(\widehat{\boldsymbol{\Sigma}}_k + \alpha \mathbf{I})^{-1}\widehat{\boldsymbol{\Sigma}}^{1/2} = (\widetilde{\boldsymbol{\Sigma}}_k^{(1)} + \alpha \widehat{\boldsymbol{\Sigma}}^{-1})^{-1} \text{ and } \boldsymbol{\varepsilon}_{t+1,k} = [\mathbf{I} - (\widetilde{\boldsymbol{\Sigma}}_k^{(1)} + \alpha \widehat{\boldsymbol{\Sigma}}^{-1})^{-1}]\boldsymbol{\varepsilon}_t.$$
  
Define  $\widehat{\mathbf{D}} = (\mathbf{I} + \alpha \widehat{\boldsymbol{\Sigma}}^{-1})^{-1} \text{ and } \widetilde{\boldsymbol{\Sigma}}_k = \widehat{\mathbf{D}}^{1/2} \widetilde{\boldsymbol{\Sigma}}_k^{(1)} \widehat{\mathbf{D}}^{1/2}.$  From

$$(\widetilde{\boldsymbol{\Sigma}}_{k}^{(1)} + \alpha \widehat{\boldsymbol{\Sigma}}^{-1})^{-1} = [\widehat{\mathbf{D}}^{-1} + (\widetilde{\boldsymbol{\Sigma}}_{k}^{(1)} - \mathbf{I})]^{-1} = \widehat{\mathbf{D}}^{1/2} [\mathbf{I} + (\widetilde{\boldsymbol{\Sigma}}_{k} - \widehat{\mathbf{D}})]^{-1} \widehat{\mathbf{D}}^{1/2}$$

and

$$\begin{split} \widetilde{\boldsymbol{\Sigma}}_k - \widehat{\mathbf{D}} &= \widehat{\mathbf{D}}^{1/2} (\widetilde{\boldsymbol{\Sigma}}_k^{(1)} - \mathbf{I}) \widehat{\mathbf{D}}^{1/2} = (\mathbf{I} + \alpha \widehat{\boldsymbol{\Sigma}}^{-1})^{-1/2} \widehat{\boldsymbol{\Sigma}}^{-1/2} (\widehat{\boldsymbol{\Sigma}}_k - \widehat{\boldsymbol{\Sigma}}) \widehat{\boldsymbol{\Sigma}}^{-1/2} (\mathbf{I} + \alpha \widehat{\boldsymbol{\Sigma}}^{-1})^{-1/2} \\ &= (\widehat{\boldsymbol{\Sigma}} + \alpha \mathbf{I})^{-1/2} (\widehat{\boldsymbol{\Sigma}}_k - \widehat{\boldsymbol{\Sigma}}) (\widehat{\boldsymbol{\Sigma}} + \alpha \mathbf{I})^{-1/2} = \boldsymbol{\Delta}_k, \end{split}$$

we get

$$\boldsymbol{\varepsilon}_{t+1,k} = (\mathbf{I} - \widehat{\mathbf{D}}^{1/2} \mathbf{C}_k \widehat{\mathbf{D}}^{1/2}) \boldsymbol{\varepsilon}_t$$
 and  $\boldsymbol{\varepsilon}_{t+1} = (\mathbf{I} - \widehat{\mathbf{D}}^{1/2} \mathbf{C} \widehat{\mathbf{D}}^{1/2}) \boldsymbol{\varepsilon}_t$ ,

where  $\mathbf{C}_k = (\mathbf{I} + \boldsymbol{\Delta}_k)^{-1}$  and  $\mathbf{C} = \frac{1}{m} \sum_{k=1}^m \mathbf{C}_k$ . Let  $\mathbf{R}_k = \mathbf{C}_k - (\mathbf{I} - \boldsymbol{\Delta}_k)$  and  $\mathbf{R} = \frac{1}{m} \sum_{k=1}^m \mathbf{R}_k = \mathbf{C} - \mathbf{I}$ . We have  $\|\mathbf{I} - \widehat{\mathbf{D}}^{1/2} \mathbf{C} \widehat{\mathbf{D}}^{1/2}\|_2 = \|\mathbf{I} - \widehat{\mathbf{D}}^{1/2} (\mathbf{I} + \mathbf{R}) \widehat{\mathbf{D}}^{1/2}\|_2$ . Below we control the right-hand side.

By  $\Delta = \max_{k \in [m]} \|\mathbf{\Delta}_k\|_2 \le 1/2$  and Lemma B.4, we obtain that  $\|\mathbf{C}_k\|_2 \le \frac{1}{1-\Delta} \le 2$  and

 $\|\mathbf{R}_k\|_2 \leq 2\Delta^2$ . Consequently,  $\|\mathbf{C}\|_2 \leq 2$  and  $\|\mathbf{R}\|_2 \leq 2\Delta^2 \leq 1/2$ . Then we obtain that

$$(1 - 2\Delta^{2})\mathbf{I} \preceq \mathbf{I} + \mathbf{R} \preceq (1 + 2\Delta^{2})\mathbf{I},$$

$$(1 - 2\Delta^{2})\widehat{\mathbf{D}} \preceq \widehat{\mathbf{D}}^{1/2}(\mathbf{I} + \mathbf{R})\widehat{\mathbf{D}}^{1/2} \preceq (1 + 2\Delta^{2})\widehat{\mathbf{D}},$$

$$\mathbf{I} - (1 + 2\Delta^{2})\widehat{\mathbf{D}} \preceq \mathbf{I} - \widehat{\mathbf{D}}^{1/2}(\mathbf{I} + \mathbf{R})\widehat{\mathbf{D}}^{1/2} \preceq \mathbf{I} - (1 - 2\Delta^{2})\widehat{\mathbf{D}}.$$

Consequently,

$$\|\mathbf{I} - \widehat{\mathbf{D}}^{1/2}(\mathbf{I} + \mathbf{R})\widehat{\mathbf{D}}^{1/2}\|_2 \le \max\left\{\|\mathbf{I} - (1 - 2\Delta^2)\widehat{\mathbf{D}}\|_2, \|\mathbf{I} - (1 + 2\Delta^2)\widehat{\mathbf{D}}\|_2\right\}.$$

Let  $\{\widehat{\lambda}_j\}_{j=1}^p$  be the eigenvalues of  $\widehat{\mathbf{\Sigma}}$  sorted in descending order. Since  $\widehat{\mathbf{D}}$  has eigenvalues  $\{(1+\alpha/\widehat{\lambda}_j)^{-1}\}_{j=1}^p\subseteq (0,1]$ , the eigenvalues of  $\mathbf{I}-(1\pm2\Delta^2)\widehat{\mathbf{D}}$  are  $\left\{1-\frac{1\pm2\Delta^2}{1+\alpha/\widehat{\lambda}_j}\right\}_{j=1}^p$ . Then

$$\|\mathbf{I} - (1 \pm 2\Delta^2)\widehat{\mathbf{D}}\|_2 = \max\left\{ \left| 1 - \frac{1 \pm 2\Delta^2}{1 + \alpha/\widehat{\lambda}_1} \right|, \left| 1 - \frac{1 \pm 2\Delta^2}{1 + \alpha/\widehat{\lambda}_p} \right| \right\}.$$

By elementary calculation and the fact  $2\Delta^2 \le 1/2 < 1$  we get

$$\begin{vmatrix} 1 - \frac{1 + 2\Delta^2}{1 + \alpha/\widehat{\lambda}_1} \end{vmatrix} = \frac{\left| \alpha/\widehat{\lambda}_1 - 2\Delta^2 \right|}{1 + \alpha/\widehat{\lambda}_1} \le \frac{\max\{\alpha/\widehat{\lambda}_1, 2\Delta^2\}}{1 + \alpha/\widehat{\lambda}_1} \le \max\left\{\frac{\alpha/\widehat{\lambda}_p}{1 + \alpha/\widehat{\lambda}_p}, \ 2\Delta^2\right\},$$

$$\begin{vmatrix} 1 - \frac{1 + 2\Delta^2}{1 + \alpha/\widehat{\lambda}_p} \end{vmatrix} = \frac{\left| \alpha/\widehat{\lambda}_p - 2\Delta^2 \right|}{1 + \alpha/\widehat{\lambda}_p} \le \frac{2\Delta^2 + \alpha/\widehat{\lambda}_p}{1 + \alpha/\widehat{\lambda}_p},$$

$$0 \le 1 - \frac{1 - 2\Delta^2}{1 + \alpha/\widehat{\lambda}_1} \le 1 - \frac{1 - 2\Delta^2}{1 + \alpha/\widehat{\lambda}_p} = \frac{2\Delta^2 + \alpha/\widehat{\lambda}_p}{1 + \alpha/\widehat{\lambda}_p}.$$

Therefore,

$$\|\mathbf{I} - \widehat{\mathbf{D}}^{1/2} \mathbf{C} \widehat{\mathbf{D}}^{1/2}\|_{2} \le \max \left\{ \frac{2\Delta^{2} + \alpha/\widehat{\lambda}_{p}}{1 + \alpha/\widehat{\lambda}_{p}}, \ 2\Delta^{2} \right\} = \frac{2\Delta^{2} + \alpha/\widehat{\lambda}_{p}}{1 + \alpha/\widehat{\lambda}_{p}} < 1.$$

**Lemma A.7.** Suppose Assumption 3.1 hold. Then there exists a constant C determined by  $\|\mathbf{u}_i\|_{\psi_2}$ , such that  $\mathbb{P}(\Delta \leq 1/2) \geq 1 - 2me^{-n/C}$  holds under either of the two conditions: (i)  $n \geq Cp$  and  $\alpha \geq 0$ ; (ii)  $\alpha \geq C\operatorname{Tr}(\mathbf{\Sigma})/n$ .

Proof of Lemma A.7. Define  $\widehat{\mathbf{S}}_k = \frac{1}{n} \sum_{i \in \mathcal{I}_k} \mathbf{u}_i \mathbf{u}_i^{\top}$  and  $\bar{\mathbf{u}}_{(k)} = \frac{1}{n} \sum_{i \in \mathcal{I}_k} \mathbf{u}_i$  for  $k \in [m]$ . Then we have  $\widehat{\boldsymbol{\Sigma}}_k = \frac{1}{n} \sum_{i \in \mathcal{I}_k} \mathbf{x}_i \mathbf{x}_i^{\top} = \begin{pmatrix} 1 & \bar{\mathbf{u}}_{(k)}^{\top} \\ \bar{\mathbf{u}}_{(k)} & \widehat{\mathbf{S}}_k \end{pmatrix}$ . Let  $\boldsymbol{\Sigma}^* = \mathbb{E}\widehat{\boldsymbol{\Sigma}}_k = \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma} \end{pmatrix}$  and observe that

$$\Delta_k = (\widehat{\Sigma} + \alpha \mathbf{I})^{-1/2} (\widehat{\Sigma}_k - \Sigma^*) (\widehat{\Sigma} + \alpha \mathbf{I})^{-1/2} - (\widehat{\Sigma} + \alpha \mathbf{I})^{-1/2} (\widehat{\Sigma} - \Sigma^*) (\widehat{\Sigma} + \alpha \mathbf{I})^{-1/2}.$$

Let  $\mathbf{B}_k = (\widehat{\Sigma} + \alpha \mathbf{I})^{-1/2} (\widehat{\Sigma}_k - \Sigma^*) (\widehat{\Sigma} + \alpha \mathbf{I})^{-1/2}$ . Since

$$\max_{k \in [m]} \|\mathbf{\Delta}_k\|_2 = \max_{k \in [m]} \|\mathbf{B}_k - \frac{1}{m} \sum_{\ell=1}^m \mathbf{B}_\ell\|_2 \le 2 \max_{k \in [m]} \|\mathbf{B}_k\|_2, \tag{A.9}$$

it boils down to bound  $\|\mathbf{B}_k\|_2$ . To do this, we let  $\mathbf{A}_0 = (\mathbf{\Sigma}^* + \alpha \mathbf{I})^{1/2} (\widehat{\mathbf{\Sigma}} + \alpha \mathbf{I})^{-1/2}$  and write

$$\|\mathbf{B}_k\|_2 = \|\mathbf{A}_0^{\top} (\mathbf{\Sigma}^* + \alpha \mathbf{I})^{-1/2} (\widehat{\mathbf{\Sigma}}_k - \mathbf{\Sigma}^*) (\mathbf{\Sigma}^* + \alpha \mathbf{I})^{-1/2} \mathbf{A}_0\|_2$$
  
$$\leq \|\mathbf{A}_0\|_2^2 \|(\mathbf{\Sigma}^* + \alpha \mathbf{I})^{-1/2} (\widehat{\mathbf{\Sigma}}_k - \mathbf{\Sigma}^*) (\mathbf{\Sigma}^* + \alpha \mathbf{I})^{-1/2} \|_2.$$

Define  $\mathbf{D} = (\mathbf{I} + \alpha(\mathbf{\Sigma}^*)^{-1})^{-1}$  and  $\widehat{\mathbf{\Sigma}}_k^{(1)} = (\mathbf{\Sigma}^* + \alpha \mathbf{I})^{-1/2} \widehat{\mathbf{\Sigma}}_k (\mathbf{\Sigma}^* + \alpha \mathbf{I})^{-1/2}$ . On the one hand,

$$(\mathbf{\Sigma}^* + \alpha \mathbf{I})^{-1/2} (\widehat{\mathbf{\Sigma}}_k - \mathbf{\Sigma}^*) (\mathbf{\Sigma}^* + \alpha \mathbf{I})^{-1/2} = \widehat{\mathbf{\Sigma}}_k^{(1)} - \mathbf{D}.$$

On the other hand,

$$\begin{split} \|\mathbf{A}_{0}\|_{2}^{2} &= \|\mathbf{A}_{0}\mathbf{A}_{0}^{\top}\|_{2} = \|(\mathbf{\Sigma}^{*} + \alpha \mathbf{I})^{1/2}(\widehat{\mathbf{\Sigma}} + \alpha \mathbf{I})^{-1}(\mathbf{\Sigma}^{*} + \alpha \mathbf{I})^{1/2}\|_{2} \\ &= \|(\mathbf{\Sigma}^{*} + \alpha \mathbf{I})^{1/2}[(\widehat{\mathbf{\Sigma}} - \mathbf{\Sigma}^{*}) + (\mathbf{\Sigma}^{*} + \alpha \mathbf{I})]^{-1}(\mathbf{\Sigma}^{*} + \alpha \mathbf{I})^{1/2}\|_{2} \\ &\leq \|[(\mathbf{\Sigma}^{*} + \alpha \mathbf{I})^{-1/2}(\widehat{\mathbf{\Sigma}} - \mathbf{\Sigma}^{*})(\mathbf{\Sigma}^{*} + \alpha \mathbf{I})^{-1/2} + \mathbf{I}]^{-1}\|_{2} = \|[\mathbf{I} + \frac{1}{m}\sum_{k=1}^{m}(\widehat{\mathbf{\Sigma}}_{k}^{(1)} - \mathbf{D})]^{-1}\|_{2} \\ &\leq \frac{1}{1 - \|\frac{1}{m}\sum_{k=1}^{m}(\widehat{\mathbf{\Sigma}}_{k}^{(1)} - \mathbf{D})\|_{2}} \leq \frac{1}{1 - \max_{k \in [m]} \|\widehat{\mathbf{\Sigma}}_{k}^{(1)} - \mathbf{D}\|_{2}}, \end{split}$$

where we used Lemma B.4. Based on these, we have

$$\max_{k \in [m]} \|\mathbf{B}_k\|_2 \le \frac{\max_{k \in [m]} \|\widehat{\boldsymbol{\Sigma}}_k^{(1)} - \mathbf{D}\|_2}{1 - \max_{k \in [m]} \|\widehat{\boldsymbol{\Sigma}}_k^{(1)} - \mathbf{D}\|_2},$$
(A.10)

and it suffices to prove under the given conditions that

$$\mathbb{P}\left(\max_{k\in[m]}\|\widehat{\boldsymbol{\Sigma}}_k^{(1)} - \mathbf{D}\|_2 \le 1/5\right) \ge 1 - 2me^{-n/C} \tag{A.11}$$

holds for some constant C.

By definition, we have

$$\widehat{\boldsymbol{\Sigma}}_{k}^{(1)} - \mathbf{D} = \begin{pmatrix} (1+\alpha)^{-1/2} & \mathbf{0} \\ \mathbf{0} & (\boldsymbol{\Sigma} + \alpha \mathbf{I})^{-1/2} \end{pmatrix} \begin{pmatrix} \mathbf{0} & \bar{\mathbf{u}}_{(k)}^{\top} \\ \bar{\mathbf{u}}_{(k)} & \widehat{\mathbf{S}}_{k} - \boldsymbol{\Sigma} \end{pmatrix} \begin{pmatrix} (1+\alpha)^{-1/2} & \mathbf{0} \\ \mathbf{0} & (\boldsymbol{\Sigma} + \alpha \mathbf{I})^{-1/2} \end{pmatrix} \\
= \begin{pmatrix} \mathbf{0} & (1+\alpha)^{-1/2} [(\boldsymbol{\Sigma} + \alpha \mathbf{I})^{-1/2} \bar{\mathbf{u}}_{(k)}]^{\top} \\ (1+\alpha)^{-1/2} (\boldsymbol{\Sigma} + \alpha \mathbf{I})^{-1/2} \bar{\mathbf{u}}_{(k)} & (\boldsymbol{\Sigma} + \alpha \mathbf{I})^{-1/2} (\widehat{\mathbf{S}}_{k} - \boldsymbol{\Sigma}) (\boldsymbol{\Sigma} + \alpha \mathbf{I})^{-1/2} \end{pmatrix}$$

and as a result,

$$\|\widehat{\mathbf{\Sigma}}_{k}^{(1)} - \mathbf{D}\|_{2} \le (1+\alpha)^{-1/2} \|(\mathbf{\Sigma} + \alpha \mathbf{I})^{-1/2} \bar{\mathbf{u}}_{(k)}\|_{2} + \|(\mathbf{\Sigma} + \alpha \mathbf{I})^{-1/2} (\widehat{\mathbf{S}}_{k} - \mathbf{\Sigma})(\mathbf{\Sigma} + \alpha \mathbf{I})^{-1/2}\|_{2}.$$
(A.12)

Here we used a simple fact that  $\left\| \begin{pmatrix} \mathbf{0} & \mathbf{A}^\top \\ \mathbf{A} & \mathbf{0} \end{pmatrix} \right\|_2 = \|\mathbf{A}\|_2$  for any matrix  $\mathbf{A}$ .

Observe that  $\mathbf{v}_i = (\mathbf{\Sigma} + \alpha \mathbf{I})^{-1/2} \mathbf{u}_i$  is a sub-gaussian random variable with zero mean and covariance matrix  $(1 + \alpha \mathbf{\Sigma}^{-1})^{-1}$ . On the other hand,  $(\mathbf{\Sigma} + \alpha \mathbf{I})^{-1/2} (\widehat{\mathbf{S}}_k - \mathbf{\Sigma}) (\mathbf{\Sigma} + \alpha \mathbf{I})^{-1/2} = \frac{1}{n} \sum_{i \in \mathcal{I}_k} \mathbf{v}_i \mathbf{v}_i^{\mathsf{T}} - (1 + \alpha \mathbf{\Sigma}^{-1})^{-1}$ . Lemma B.5 forces

$$\mathbb{P}\left(\|(\mathbf{\Sigma} + \alpha \mathbf{I})^{-1/2} \bar{\mathbf{u}}_{(k)}\|_{2} > 1/10\right) \leq e^{-n/C},$$

$$\mathbb{P}\left(\|(\mathbf{\Sigma} + \alpha \mathbf{I})^{-1/2} (\widehat{\mathbf{S}}_{k} - \mathbf{\Sigma})(\mathbf{\Sigma} + \alpha \mathbf{I})^{-1/2}\|_{2} > 1/10\right) \leq e^{-n/C},$$

where C is the constant therein. These estimates and (A.12) lead to (A.11).

Following the similar idea in the proof above, we get the following results.

**Lemma A.8.** Suppose Assumption 3.1 holds with C being the constant in Lemma A.7. Then

- $\mathbb{P}(\|(\mathbf{\Sigma}^*)^{-1/2}(\widehat{\mathbf{\Sigma}} \mathbf{\Sigma}^*)(\mathbf{\Sigma}^*)^{-1/2}\|_2 \le 1/2) \ge 1 2me^{-N/C};$
- $\mathbb{P}(\Delta \leq C_2'\sqrt{p/n}) \geq 1 2e^{-C_1'p}$  holds for some constants  $C_1'$  and  $C_2'$

Now we come back to the main proof. We will use the three lemmas above to show that

with high probability,  $\lambda_{\max}(\widehat{\Sigma})/\lambda_{\min}(\widehat{\Sigma}) \leq 3\kappa$  and

$$\left\|\widehat{\mathbf{\Sigma}}^{1/2}\left(\boldsymbol{\theta}_{t+1} - \widehat{\boldsymbol{\theta}}\right)\right\|_{2} \le \left(1 - \frac{1 - \min\{1/2, C_{2}p/n\}}{1 + C_{3}\alpha}\right) \left\|\widehat{\mathbf{\Sigma}}^{1/2}\left(\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\right)\right\|_{2}, \quad \forall t \ge 0. \quad (A.13)$$

Then we conclude the proof by induction and some simple linear algebra.

First, Lemma A.8 asserts that with high probability,  $\|(\mathbf{\Sigma}^*)^{-1/2}(\widehat{\mathbf{\Sigma}} - \mathbf{\Sigma}^*)(\mathbf{\Sigma}^*)^{-1/2}\|_2 \le 1/2$ . On this event, we have  $-\frac{1}{2}\mathbf{\Sigma}^* \le \widehat{\mathbf{\Sigma}} - \mathbf{\Sigma}^* \le \frac{1}{2}\mathbf{\Sigma}^*$  and thus  $\frac{1}{2}\mathbf{\Sigma}^* \le \widehat{\mathbf{\Sigma}} \le \frac{3}{2}\mathbf{\Sigma}^*$ . Hence

$$\frac{1}{2}\lambda_{\min}(\mathbf{\Sigma}^*) \leq \lambda_{\min}(\widehat{\mathbf{\Sigma}}) \leq \lambda_{\max}(\widehat{\mathbf{\Sigma}}) \leq \frac{3}{2}\lambda_{\max}(\mathbf{\Sigma}^*)$$

and  $\lambda_{\max}(\widehat{\Sigma})/\lambda_{\min}(\widehat{\Sigma}) \leq 3\lambda_{\max}(\Sigma^*)/\lambda_{\min}(\Sigma^*) = 3\kappa$ .

Second, Lemma A.6 forces

$$\left\|\widehat{\Sigma}^{1/2} \left(\boldsymbol{\theta}_{t+1} - \widehat{\boldsymbol{\theta}}\right)\right\|_{2} \leq \frac{2\Delta^{2} + \alpha/\lambda_{\min}(\widehat{\boldsymbol{\Sigma}})}{1 + \alpha/\lambda_{\min}(\widehat{\boldsymbol{\Sigma}})} \left\|\widehat{\boldsymbol{\Sigma}}^{1/2} \left(\boldsymbol{\theta}_{t} - \widehat{\boldsymbol{\theta}}\right)\right\|_{2}, \quad \forall t \geq 0.$$
 (A.14)

Lemmas A.7 and A.8 imply that  $\Delta \leq 1/2$ ,  $\Delta \leq \tilde{C}\sqrt{p/n}$ , and  $\lambda_{\min}(\widehat{\Sigma}) \geq 1$  hold simultaneously with high probability, where  $\tilde{C}$  is some constant. On this event, we have

$$\frac{2\Delta^2 + \alpha/\lambda_{\min}(\widehat{\Sigma})}{1 + \alpha/\lambda_{\min}(\widehat{\Sigma})} = 1 - \frac{1 - 2\Delta^2}{1 + \alpha/\lambda_{\min}(\widehat{\Sigma})} \le 1 - \frac{1 - \min\{1/2, C_2 p/n\}}{1 + C_3 \alpha}$$

for some constants  $C_2$  and  $C_3$ . Then we get (A.13) from the estimates above and complete the proof.

## B Technical lemmas

The following lemma lists basic properties of strongly convex functions, which can be found in standard textbooks on convex optimization (Nesterov, 2013).

**Lemma B.1.** Suppose f is a convex function defined on some convex open set  $\Omega \subseteq \mathbb{R}^p$ , and  $\partial f(\mathbf{x})$  denotes its subdifferential set at  $\mathbf{x} \in \Omega$ . The followings are equivalent:

- f is  $\rho$ -strongly convex in  $\Omega$ ;
- $f[(1-t)\mathbf{x} + t\mathbf{y}] \le (1-t)f(\mathbf{x}) + tf(\mathbf{y}) (\rho/2)t(1-t)\|\mathbf{y} \mathbf{x}\|_2^2$ ,  $\forall \mathbf{x}, \mathbf{y} \in \Omega$  and  $t \in [0, 1]$ ;
- $\langle \mathbf{h} \mathbf{g}, \mathbf{y} \mathbf{x} \rangle \ge \rho \|\mathbf{y} \mathbf{x}\|_2^2$ ,  $\forall \mathbf{x}, \mathbf{y} \in \Omega$ ,  $\mathbf{g} \in \partial f(\mathbf{x})$  and  $\mathbf{h} \in \partial f(\mathbf{y})$ .

If any of the above holds, f is said to be  $\rho$ -strongly convex.

**Lemma B.2.** Let  $f: \mathbb{R}^p \to \mathbb{R}$  be a convex function. Suppose there exists  $\mathbf{x} \in \mathbb{R}^p$  and r > 0 such that f is  $\rho$ -strongly convex in  $B(\mathbf{x}, r)$ . If  $\|\mathbf{h} - \mathbf{g}\|_2 < \rho r$  holds for some  $\mathbf{g} \in \partial f(\mathbf{x})$  and  $\mathbf{h} \in \partial f(\mathbf{y})$ , then  $\|\mathbf{y} - \mathbf{x}\|_2 \leq \|\mathbf{h} - \mathbf{g}\|_2/\rho \leq r$ .

*Proof.* If we know a priori that  $\|\mathbf{y} - \mathbf{x}\|_2 \le r$ , then we use the strong convexity of f in  $B(\mathbf{x}, r)$  and Cauchy-Schwarz inequality to obtain

$$\rho \|\mathbf{y} - \mathbf{x}\|_2^2 \le \langle \mathbf{h} - \mathbf{g}, \mathbf{y} - \mathbf{x} \rangle \le \|\mathbf{h} - \mathbf{g}\|_2 \|\mathbf{y} - \mathbf{x}\|_2,$$

and get the desired result. Suppose on the contrary that  $\|\mathbf{y} - \mathbf{x}\|_2 > r$ , and define  $\bar{\mathbf{y}} = \mathbf{x} + r(\mathbf{y} - \mathbf{x})/\|\mathbf{y} - \mathbf{x}\|_2$ . Then  $\|\bar{\mathbf{y}} - \mathbf{x}\|_2 = r$ . The strong convexity of f in  $B(\mathbf{x}, r)$  and Lemma B.1 yield

$$\langle \mathbf{s} - \mathbf{g}, \bar{\mathbf{y}} - \mathbf{x} \rangle \ge \rho \|\bar{\mathbf{y}} - \mathbf{x}\|_2^2, \quad \forall \mathbf{s} \in \partial f(\bar{\mathbf{y}}).$$

By the convexity of f, we always have

$$\langle \mathbf{h} - \mathbf{s}, \bar{\mathbf{y}} - \mathbf{x} \rangle = \frac{r}{\|\mathbf{y} - \mathbf{x}\|_2 - r} \langle \mathbf{h} - \mathbf{s}, \mathbf{y} - \bar{\mathbf{y}} \rangle \ge 0, \quad \forall \mathbf{s} \in \partial f(\bar{\mathbf{y}}).$$

Summing up the two inequalities above, we get

$$\rho \|\bar{\mathbf{y}} - \mathbf{x}\|_2^2 \le \langle \mathbf{h} - \mathbf{g}, \bar{\mathbf{y}} - \mathbf{x} \rangle \le \|\mathbf{h} - \mathbf{g}\|_2 \|\bar{\mathbf{y}} - \mathbf{x}\|_2,$$

where we also used the Cauchy-Schwarz inequality. Then  $\|\mathbf{h} - \mathbf{g}\|_2 \ge \rho \|\bar{\mathbf{y}} - \mathbf{x}\|_2 = \rho r$  leads to contradiction. Hence, we must have only the case  $\|\mathbf{y} - \mathbf{x}\|_2 \le r$ .

**Lemma B.3.** Let  $f: \mathbb{R}^p \to \mathbb{R}$  be a convex function. For any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$ , we have

$$\|\operatorname{prox}_f(\mathbf{x}) - \operatorname{prox}_f(\mathbf{y})\|_2^2 \le \langle \mathbf{x} - \mathbf{y}, \operatorname{prox}_f(\mathbf{x}) - \operatorname{prox}_f(\mathbf{y}) \rangle.$$

If  $\inf_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x}) > -\infty$ , then  $\|\operatorname{prox}_f(\mathbf{x}) - \mathbf{x}^*\|_2 \le \|\mathbf{x} - \mathbf{x}^*\|_2$  and

$$\|\operatorname{prox}_{f}(\mathbf{x}) - \mathbf{x}\|_{2}^{2} \leq \|\mathbf{x} - \mathbf{x}^{*}\|_{2}^{2} - \|\operatorname{prox}_{f}(\mathbf{x}) - \mathbf{x}^{*}\|_{2}^{2}$$

hold for any  $\mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x})$ .

If f is  $\rho$ -strongly convex in  $B(\mathbf{x}^*, r)$  for some r > 0 and  $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x})$ , then  $\|\operatorname{prox}_{\alpha^{-1}f}(\mathbf{x}) - \mathbf{x}^*\|_2 \le \frac{\alpha}{\alpha + \rho} \|\mathbf{x} - \mathbf{x}^*\|_2$ ,  $\forall \mathbf{x} \in B(\mathbf{x}^*, r)$  and  $\alpha > 0$ .

**Proof of Lemma B.3.** The first claim is the well-known "firm non-expansiveness" property of the proximal mapping (Parikh and Boyd, 2014).

If  $\inf_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x}) > -\infty$ , then any  $\mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x})$  is a fixed point of  $\operatorname{prox}_f$ . The firm non-expansiveness with  $\mathbf{y} = \mathbf{x}^*$  yields

$$\|\operatorname{prox}_{f}(\mathbf{x}) - \mathbf{x}^{*}\|_{2}^{2} \le \langle \mathbf{x} - \mathbf{x}^{*}, \operatorname{prox}_{f}(\mathbf{x}) - \mathbf{x}^{*} \rangle \tag{B.1}$$

and then  $\|\operatorname{prox}_f(\mathbf{x}) - \mathbf{x}^*\|_2 \le \|\mathbf{x} - \mathbf{x}^*\|_2$ . The next claim is proved by

$$\begin{aligned} \| \operatorname{prox}_{f}(\mathbf{x}) - \mathbf{x} \|_{2}^{2} &= \| [\operatorname{prox}_{f}(\mathbf{x}) - \mathbf{x}^{*}] - (\mathbf{x} - \mathbf{x}^{*}) \|_{2}^{2} \\ &= \| \operatorname{prox}_{f}(\mathbf{x}) - \mathbf{x}^{*} \|_{2}^{2} + \| \mathbf{x} - \mathbf{x}^{*} \|_{2}^{2} - 2 \langle \operatorname{prox}_{f}(\mathbf{x}) - \mathbf{x}^{*}, \mathbf{x} - \mathbf{x}^{*} \rangle \\ &\leq \| \operatorname{prox}_{f}(\mathbf{x}) - \mathbf{x}^{*} \|_{2}^{2} + \| \mathbf{x} - \mathbf{x}^{*} \|_{2}^{2} - 2 \| \operatorname{prox}_{f}(\mathbf{x}) - \mathbf{x}^{*} \|_{2}^{2} \\ &= \| \mathbf{x} - \mathbf{x}^{*} \|_{2}^{2} - \| \operatorname{prox}_{f}(\mathbf{x}) - \mathbf{x}^{*} \|_{2}^{2}, \end{aligned}$$

where the inequality follows from (B.1).

For the last claim, we fix any  $\alpha > 0$  and  $\mathbf{x} \in B(\mathbf{x}^*, r)$  and define  $\mathbf{x}^+ = \operatorname{prox}_{\alpha^{-1}f}(\mathbf{x})$ . Then  $\|\mathbf{x}^+ - \mathbf{x}^*\|_2 \le \|\mathbf{x} - \mathbf{x}^*\|_2 < r$ . The optimality conditions for  $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{y} \in \mathbb{R}^p} f(\mathbf{y})$  and  $\mathbf{x}^+ = \operatorname{argmin}_{\mathbf{y} \in \mathbb{R}^p} \{f(\mathbf{y}) + (\alpha/2) \|\mathbf{y} - \mathbf{x}\|_2^2\}$  imply that  $\mathbf{0} \in \partial f(\mathbf{x}^*)$  and  $-\alpha(\mathbf{x}^+ - \mathbf{x}) \in \partial f(\mathbf{x}^+)$ . Since f is  $\rho$ -strongly convex in  $B(\mathbf{x}^*, r)$ , Lemma B.1 forces

$$\rho \|\mathbf{x}^{+} - \mathbf{x}^{*}\|_{2}^{2} \leq \langle -\alpha(\mathbf{x}^{+} - \mathbf{x}) - \mathbf{0}, \mathbf{x}^{+} - \mathbf{x}^{*} \rangle = -\alpha \|\mathbf{x}^{+} - \mathbf{x}^{*}\|_{2}^{2} - \alpha \langle \mathbf{x}^{*} - \mathbf{x}, \mathbf{x}^{+} - \mathbf{x}^{*} \rangle$$
$$\leq -\alpha \|\mathbf{x}^{+} - \mathbf{x}^{*}\|_{2}^{2} + \alpha \|\mathbf{x}^{*} - \mathbf{x}\|_{2} \|\mathbf{x}^{+} - \mathbf{x}^{*}\|_{2}$$

and thus  $\|\mathbf{x}^+ - \mathbf{x}^*\|_2 \le \frac{\alpha}{\alpha + \rho} \|\mathbf{x} - \mathbf{x}^*\|_2$ .

**Lemma B.4** (Neumann expansion). Let  $\|\cdot\|$  be a submultiplicative matrix norm with  $\|\mathbf{I}\| = 1$ . When  $\|\mathbf{M}\| < 1$ , we have  $(\mathbf{I} - \mathbf{M})^{-1} = \sum_{j=0}^{\infty} \mathbf{M}^j = \mathbf{I} + \mathbf{M} + \mathbf{M}(\mathbf{I} - \mathbf{M})^{-1}\mathbf{M}$ ,  $\|\mathbf{I} - \mathbf{M}\| \le 1/(1 - \|\mathbf{M}\|)$  and  $\|(\mathbf{I} - \mathbf{M})^{-1} - (\mathbf{I} + \mathbf{M})\| \le \|\mathbf{M}\|^2/(1 - \|\mathbf{M}\|)$ .

**Lemma B.5.** Let  $\mathbf{S} \succ 0$  and  $\alpha \geq 0$  be deterministic,  $\mathbf{A} = (\mathbf{I} + \alpha \mathbf{S}^{-1})^{-1}$ ,  $\{\mathbf{u}_i\}_{i=1}^n \subseteq \mathbb{R}^d$  be i.i.d. sub-gaussian random vectors with zero mean and covariance matrix  $\mathbf{A}$ ,  $\bar{\mathbf{u}} = \frac{1}{n} \sum_{i=1}^n \mathbf{u}_i$  and  $\widehat{\mathbf{A}} = \frac{1}{n} \sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^T$ . Then the following hold:

- 1. There exists some positive constant C that only depends on  $\|\mathbf{u}_1\|_{\psi_2}$ , such that  $\mathbb{P}(\|\bar{\mathbf{u}}\|_2 > 1/10) \leq e^{-n/C}$  and  $\mathbb{P}(\|\hat{\mathbf{A}} \mathbf{A}\|_2 > 1/10) \leq e^{-n/C}$  hold under any one of the two conditions holds: (i)  $n \geq Cd$ ; (ii)  $\alpha \geq C \operatorname{Tr}(\mathbf{S})/n$ .
- 2. Suppose that  $n \geq cd$  for some constant c > 0. There exist positive constants  $C_1'$  and  $C_2'$  such that  $\mathbb{P}(\max\{\|\bar{\mathbf{u}}\|_2, \|\widehat{\mathbf{A}} \mathbf{A}\|_2\} \geq C_2' \sqrt{d/n}) \leq 2e^{-C_1'd}$ .

**Proof of Lemma B.5.** Let  $\lambda_1 > \cdots > \lambda_d > 0$  be the eigenvalues of **S**, and C > 0 be some constant to be determined. When  $n \geq Cd$ , the fact  $\|\mathbf{A}\|_2 = 1/(1 + \alpha/\lambda_1) \leq 1$  implies  $n \geq C \text{Tr}(\mathbf{A})$ . Also, if  $\alpha \geq C \text{Tr}(\mathbf{S})/n$ , then the crude estimate

$$\operatorname{Tr}(\mathbf{A}) = \sum_{j=1}^{d} \frac{1}{1 + \alpha/\lambda_j} \le \sum_{j=1}^{d} \frac{1}{\alpha/\lambda_j} = \sum_{j=1}^{d} \frac{\lambda_j}{\alpha} = \frac{\operatorname{Tr}(\mathbf{S})}{\alpha} \le \frac{n}{C_0}$$

also leads to  $n \geq C\operatorname{Tr}(\mathbf{A})$ . Hence it suffices to find some proper C and show the desired results given  $n \geq C\operatorname{Tr}(\mathbf{A})$ .

Now we prove the first statement. We first study concentration of the sample mean vector  $\bar{\mathbf{u}}$ . Since  $\bar{\mathbf{u}}$  is a sub-gaussian random vector with covariance matrix  $n^{-1}\mathbf{A}$ , Theorem 2.1 in Hsu et al. (2012) asserts the existence of a constant  $c_1 > 0$  such that

$$\mathbb{P}\left[\|\bar{\mathbf{u}}\|_{2}^{2} \le c_{1}n^{-1}\left(\text{Tr}(\mathbf{A}) + 2\sqrt{\text{Tr}(\mathbf{A}^{2})t} + 2\|\mathbf{A}\|_{2}t\right)\right] \ge 1 - e^{-t}, \quad \forall t > 0.$$
 (B.2)

Choose any constant  $C_1 \geq 500c_1$ . Let  $t = n/C_1$ , and suppose that  $n \geq C_1 \text{Tr}(\mathbf{A})$ . Using  $\|\mathbf{A}\|_2 \leq 1$  and  $\text{Tr}(\mathbf{A}^2) \leq \text{Tr}(\mathbf{A}) \|\mathbf{A}\|_2 \leq \text{Tr}(\mathbf{A})$ , we get

$$c_{1}n^{-1}\left(\operatorname{Tr}(\mathbf{A}) + 2\sqrt{\operatorname{Tr}(\mathbf{A}^{2})t} + 2\|\mathbf{A}\|_{2}t\right) \leq \frac{c_{1}}{C_{1}} + \frac{2c_{1}\sqrt{\operatorname{Tr}(\mathbf{A})t}}{n} + \frac{2c_{1}t}{n}$$

$$= \frac{c_{1}}{C_{1}} + \frac{2c_{1}\sqrt{\operatorname{Tr}(\mathbf{A})n/C_{1}}}{n} + \frac{2c_{1}(n/C_{1})}{n} = \frac{3c_{1}}{C_{1}} + 2c_{1}\sqrt{\frac{\operatorname{Tr}(\mathbf{A})}{C_{1}n}} \leq \frac{5c_{1}}{C_{1}} \leq \frac{1}{10^{2}}.$$

Hence  $\mathbb{P}(\|\bar{\mathbf{u}}\|_2 > 1/10) \ge e^{-n/C_1}$ .

Now we come to concentration of the sample covariance matrix  $\widehat{\mathbf{A}}$ . Let  $r(\mathbf{A}) = \text{Tr}(\mathbf{A})/\|\mathbf{A}\|_2$ . According to the Theorem 9 in Koltchinskii and Lounici (2014), there exists a constant  $c_2 \geq 1$  such that the following holds: for any  $t \geq 1$ , with probability at least  $1 - e^{-t}$  we have

$$\|\widehat{\mathbf{A}} - \mathbf{A}\|_{2} \le c_{2} \|\mathbf{A}\|_{2} \max \left\{ \sqrt{\frac{r(\mathbf{A})}{n}}, \frac{r(\mathbf{A})}{n}, \sqrt{\frac{t}{n}}, \frac{t}{n} \right\}.$$
 (B.3)

Note that the upper bound above can be rewritten as

$$c_{2} \max \left\{ \sqrt{\frac{\|\mathbf{A}\|_{2} \operatorname{Tr}(\mathbf{A})}{n}}, \frac{\operatorname{Tr}(\mathbf{A})}{n}, \|\mathbf{A}\|_{2} \sqrt{\frac{t}{n}}, \|\mathbf{A}\|_{2} \frac{t}{n} \right\} \leq c_{2} \max \left\{ \sqrt{\frac{\operatorname{Tr}(\mathbf{A})}{n}}, \frac{\operatorname{Tr}(\mathbf{A})}{n}, \sqrt{\frac{t}{n}}, \frac{t}{n} \right\}.$$
(B.4)

Let  $C_2 = 100c_2^2$ . When  $n \ge C_2 \text{Tr}(\mathbf{A})$ , by taking  $t = n/C_2$  we get  $\mathbb{P}(\|\widehat{\mathbf{A}} - \mathbf{A}\|_2 > 1/10) \le e^{-n/C_2}$ . The proof of the first statement is then finished by taking  $C = \max\{C_1, C_2\}$ .

We proceed to prove the second statement. Let  $t = C_1'd$  for some constant  $C_1'$ . Note that  $d \geq \text{Tr}(\mathbf{A})$  and  $t \geq C_1'\text{Tr}(\mathbf{A})$ . According to (B.2), (B.3) and (B.4), we obtain that with probability at least  $1 - 2e^{-C_1'd}$ ,  $\|\bar{\mathbf{u}}\|_2 \leq \tilde{C}\sqrt{d/n}$  and  $\|\hat{\mathbf{A}} - \mathbf{A}\|_2 \leq \tilde{C} \max\left\{\sqrt{d/n}, d/n\right\}$  hold with some constant  $\tilde{C}$ . Since  $n \geq cd$ , we have

$$\max\left\{\sqrt{d/n}, d/n\right\} \le \max\{1, 1/\sqrt{c}\}\sqrt{d/n}.$$

By combining the inequalities above, we obtain that

$$\mathbb{P}(\max{\{\|\bar{\mathbf{u}}\|_2, \|\hat{\mathbf{A}} - \mathbf{A}\|_2\}} \ge C_2' \sqrt{d/n}) \le 2e^{-C_1'd}$$

with 
$$C_2' = \tilde{C} \max\{1, 1/\sqrt{c}\}.$$

## References

Arjevani, Y. and Shamir, O. (2015). Communication complexity of distributed convex learning and optimization. In *Advances in neural information processing systems*.

Banerjee, M., Durot, C., Sen, B. et al. (2019). Divide and conquer in nonstandard problems and the super-efficiency phenomenon. *The Annals of Statistics* 47 720–757.

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

Bertsekas, D. P. and Tsitsiklis, J. N. (1989). Parallel and distributed computation: numerical methods, vol. 23. Prentice hall Englewood Cliffs, NJ.

- BICKEL, P. J. (1975). One-step huber estimates in the linear model. *Journal of the American Statistical Association* **70** 428–434.
- BOYD, S., PARIKH, N., CHU, E., PELEATO, B. and ECKSTEIN, J. (2011). Distributed optimization and statistical learning via the alternating direction method of multipliers. Foundations and Trends® in Machine learning 3 1–122.
- CHEN, X., LIU, W. and ZHANG, Y. (2018). First-order newton-type estimator for distributed estimation and inference. arXiv preprint arXiv:1811.11368.
- Chen, X. and Xie, M.-g. (2014). A split-and-conquer approach for analysis of extraordinarily large data. *Statistica Sinica* 1655–1684.
- CHEN, Y., Su, L. and Xu, J. (2017). Distributed statistical machine learning in adversarial settings: Byzantine gradient descent. *Proceedings of the ACM on Measurement and Analysis of Computing Systems* 1 44.
- Crane, R. and Roosta, F. (2019). Dingo: Distributed newton-type method for gradient-norm optimization. arXiv preprint arXiv:1901.05134.
- Deng, W. and Yin, W. (2016). On the global and linear convergence of the generalized alternating direction method of multipliers. *Journal of Scientific Computing* **66** 889–916.
- Dua, D. and Graff, C. (2017). UCI machine learning repository. URL http://archive.ics.uci.edu/ml
- Duchi, J. C., Agarwal, A. and Wainwright, M. J. (2012). Dual averaging for distributed optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic control* **57** 592–606.
- FAN, J. and Li, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association* **96** 1348–1360.
- FAN, J., WANG, D., WANG, K. and ZHU, Z. (2017). Distributed estimation of principal eigenspaces. arXiv preprint arXiv:1702.06488.
- Garber, D., Shamir, O. and Srebro, N. (2017). Communication-efficient algorithms for distributed stochastic principal component analysis. arXiv preprint arXiv:1702.08169

- HAN, Y., MUKHERJEE, P., OZGUR, A. and WEISSMAN, T. (2018). Distributed statistical estimation of high-dimensional and nonparametric distributions. In 2018 IEEE International Symposium on Information Theory (ISIT). IEEE.
- HOERL, A. E. and KENNARD, R. W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics* **12** 55–67.
- Hong, M. and Luo, Z.-Q. (2017). On the linear convergence of the alternating direction method of multipliers. *Mathematical Programming* **162** 165–199.
- HSU, D., KAKADE, S. and ZHANG, T. (2012). A tail inequality for quadratic forms of subgaussian random vectors. *Electronic Communications in Probability* 17.
- JAGGI, M., SMITH, V., TAKÁC, M., TERHORST, J., KRISHNAN, S., HOFMANN, T. and JORDAN, M. I. (2014). Communication-efficient distributed dual coordinate ascent. In Advances in neural information processing systems.
- JORDAN, M. I., LEE, J. D. and YANG, Y. (2018). Communication-efficient distributed statistical inference. *Journal of the American Statistical Association*.
- KLEINER, A., TALWALKAR, A., SARKAR, P. and JORDAN, M. I. (2014). A scalable bootstrap for massive data. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **76** 795–816.
- Koltchinskii, V. and Lounici, K. (2014). Concentration inequalities and moment bounds for sample covariance operators. arXiv preprint arXiv:1405.2468.
- Konečný, J., McMahan, B. and Ramage, D. (2015). Federated optimization: Distributed optimization beyond the datacenter. arXiv preprint arXiv:1511.03575.
- Lan, G., Lee, S. and Zhou, Y. (2018). Communication-efficient algorithms for decentralized and stochastic optimization. *Mathematical Programming* 1–48.
- LEE, J. D., LIN, Q., MA, T. and YANG, T. (2017a). Distributed stochastic variance reduced gradient methods by sampling extra data with replacement. *The Journal of Machine Learning Research* **18** 4404–4446.
- LEE, J. D., LIU, Q., SUN, Y. and TAYLOR, J. E. (2017b). Communication-efficient sparse regression. *Journal of Machine Learning Research* 18 1–30.

- Mahajan, D., Agrawal, N., Keerthi, S. S., Sellamanickam, S. and Bottou, L. (2015). An efficient distributed learning algorithm based on effective local functional approximations. *Journal of Machine Learning Research* **16** 1–36.
- Mei, S., Bai, Y. and Montanari, A. (2018). The landscape of empirical risk for nonconvex losses. *The Annals of Statistics* **46** 2747–2774.
- NESTEROV, Y. (2013). Introductory lectures on convex optimization: A basic course, vol. 87. Springer Science & Business Media.
- NESTEROV, Y. E. (1983). A method for solving the convex programming problem with convergence rate o (1/k<sup>2</sup>). In *Dokl. Akad. Nauk SSSR*, vol. 269.
- NOCEDAL, J. and Wright, S. J. (2006). Numerical optimization 2nd.
- Parikh, N. and Boyd, S. (2014). Proximal algorithms. Foundations and Trends® in Optimization 1 127–239.
- RECHT, B., RE, C., WRIGHT, S. and NIU, F. (2011). Hogwild: A lock-free approach to parallelizing stochastic gradient descent. In *Advances in neural information processing systems*.
- RICHTÁRIK, P. and TAKÁČ, M. (2016). Distributed coordinate descent method for learning with big data. The Journal of Machine Learning Research 17 2657–2681.
- ROBINSON, P. M. (1988). The stochastic difference between econometric statistics. *Econometrica: Journal of the Econometric Society* 531–548.
- ROCKAFELLAR, R. T. (1976). Monotone operators and the proximal point algorithm. SIAM journal on control and optimization 14 877–898.
- ROSENBLATT, J. D. and NADLER, B. (2016). On the optimality of averaging in distributed statistical learning. *Information and Inference: A Journal of the IMA* 5 379–404.
- Shamir, O., Srebro, N. and Zhang, T. (2014). Communication-efficient distributed optimization using an approximate Newton-type method. In *International conference on machine learning*.
- Shang, Z. and Cheng, G. (2017). Computational limits of a distributed algorithm for smoothing spline. The Journal of Machine Learning Research 18 3809–3845.

- SHI, C., Lu, W. and Song, R. (2018). A massive data framework for m-estimators with cubic-rate. *Journal of the American Statistical Association* 1–12.
- SMITH, V., FORTE, S., CHENXIN, M., TAKÁČ, M., JORDAN, M. I. and JAGGI, M. (2018). Cocoa: A general framework for communication-efficient distributed optimization. *Journal of Machine Learning Research* 18 230.
- Suchard, M. A., Wang, Q., Chan, C., Frelinger, J., Cron, A. and West, M. (2010). Understanding gpu programming for statistical computation: Studies in massively parallel massive mixtures. *Journal of computational and graphical statistics* **19** 419–438.
- SZABO, B. and VAN ZANTEN, H. (2017). An asymptotic analysis of distributed nonparametric methods. arXiv preprint arXiv:1711.03149.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* 267–288.
- Volgushev, S., Chao, S.-K. and Cheng, G. (2017). Distributed inference for quantile regression processes. arXiv preprint arXiv:1701.06088.
- Wang, J., Kolar, M., Srebro, N. and Zhang, T. (2017a). Efficient distributed learning with sparsity. In *International Conference on Machine Learning*.
- Wang, J., Wang, W. and Srebro, N. (2017b). Memory and communication efficient distributed stochastic optimization with minibatch prox. In *Conference on Learning Theory*.
- Wang, S., Roosta-Khorasani, F., Xu, P. and Mahoney, M. W. (2018a). Giant: Globally improved approximate newton method for distributed optimization. In *Advances in Neural Information Processing Systems*.
- Wang, X. and Dunson, D. B. (2013). Parallelizing mcmc via weierstrass sampler. arXiv preprint arXiv:1312.4605.
- Wang, X., Yang, Z., Chen, X. and Liu, W. (2018b). Distributed inference for linear support vector machine. arXiv preprint arXiv:1811.11922.
- Woodworth, B., Wang, J., McMahan, B. and Srebro, N. (2018). Graph oracle models, lower bounds, and gaps for parallel stochastic optimization. arXiv preprint arXiv:1805.10222.

- YANG, T. (2013). Trading computation for communication: Distributed stochastic dual coordinate ascent. In Advances in Neural Information Processing Systems.
- Zhang, C.-H. et al. (2010). Nearly unbiased variable selection under minimax concave penalty. The Annals of statistics 38 894–942.
- Zhang, Y., Duchi, J. and Wainwright, M. (2015). Divide and conquer kernel ridge regression: A distributed algorithm with minimax optimal rates. *The Journal of Machine Learning Research* **16** 3299–3340.
- Zhang, Y., Duchi, J. C. and Wainwright, M. J. (2013). Communication-efficient algorithms for statistical optimization. *Journal of Machine Learning Research* **14** 3321–3363.
- Zhang, Y. and Xiao, L. (2015). Disco: Distributed optimization for self-concordant empirical loss. In *International conference on machine learning*.