

Distributed Sparse Optimization Based on Minimax Concave and Consensus Promoting Penalties: Towards Global Optimality

Kei Komuro Masahiro Yukawa

Department of Electronics and Electrical Engineering
Keio University, Japan

Renato L.G. Cavalcante

Fraunhofer Institute for Telecommunications
Heinrich Hertz Institute, Germany

Abstract—We propose a distributed optimization framework to generate accurate sparse estimates while allowing an algorithmic solution with guaranteed convergence to a global minimizer. To this end, the proposed problem formulation involves the minimax concave penalty together with an additional penalty called *consensus promoting penalty (CPP)* that induces convexity to the resulting optimization problem. This problem is solved with an exact first-order proximal gradient algorithm, which employs a pair of proximity operators and is referred to as the distributed proximal and debiasing-gradient (DPD) method. Numerical examples show that CPP not only convexifies the whole cost function, but it also accelerates the convergence speed with respect to the system mismatch.

Index Terms—distributed optimization, nonconvex penalty, Moreau envelope, proximity operator, sparseness

I. INTRODUCTION

Distributed optimization [2], [3] is a key component in many applications such as big data analytics [4] and sensor networks [5]. In some scenarios, such as environmental modeling, the estimand can be assumed sparse, so convex penalties have been used in optimization problems to promote sparsity of the solutions. Such convex formulations can be addressed with algorithms such as the proximal distributed gradient descent (Prox-DGD) [6] and the proximal gradient exact first-order algorithm (PG-EXTRA) [7], to name a few. However, despite the tractability of the convergence analysis, convex penalties typically increase the estimation bias to reduce the variance of the estimates.

In the context of nondistributed optimization, nonconvex penalties [8], [9] have widely been studied to reduce estimation biases. *Weakly-convex* penalties [10]–[13], in particular, allow algorithmic solutions with global optimality because the convexity of the overall cost can be preserved if the loss function is strongly convex. In typical distributed settings, however, the amount of data available at each node is limited, so the local loss function is not necessarily strongly convex, and hence each local cost may not be convex.

To address the above challenge, we can use approaches dealing with nonsmooth nonconvex penalties. Most, if not all, of the previous studies on nonconvex distributed approaches [14]–[16] have shown convergence to a stationary

This work was supported by JST SICORP Grant Number JPMJSC20C6, Japan. The authors also acknowledge the financial support by the Federal Ministry of Education and Research of Germany (BMBF) under grant 01DR21009 and the program “Souverän. Digital. Vernetzt.” Joint project 6G-RIC, project identification number: 16KISK020K.

A full version of this work is given in [1].

point. Examples of approaches of this type include the Prox-DGD [17] and the in-network nonconvex optimization (NEXT) [18]. Although the approaches in [19] and [20] have shown convergence to a second-order optimal point, they require smoothness of the cost functions.

In our recent work [21], the minimax concave (MC) penalty [10], [11] has been used based on the following idea. An application of Moreau’s decomposition decouples the MC penalty into a sum of a convex function and a negative squared ℓ_2 norm. The idea therein was to approximate this concave quadratic term by the sum of the squared inner product between the estimate and each input vector, where the approximation relies on the statistical orthogonality assumption of input vectors. The resultant penalty function, which is referred to as the approximate MC (AMC) penalty, convexifies the local cost, thereby ensuring convergence to a global minimizer. Despite global optimality, AMC tends to yield less accurate estimates than the original MC penalty particularly when the underlying assumption is violated. This fact motivates us to devise a better formulation using the MC penalty while guaranteeing global optimality.

In this paper, we propose a different formulation that utilizes the MC penalty without the need for approximations. Specifically, we introduce an additional term which we call *the consensus promoting penalty (CPP)* that makes each local cost function convex *with respect to each local variable*. In addition, we show by simulations that CPP will make the whole “global” cost function convex. The proposed formulation is referred to as the MC-CPP formulation. Numerical examples suggest that the MC-CPP cost function is convex, and the estimates generated by the PG-EXTRA algorithm implementing MC-CPP achieve the same level of errors as the centralized solution at the steady state.

II. PRELIMINARIES

This section introduces notation, convex analytic tools, and the distributed sparse optimization problem.

A. Notation

Let \mathbb{R}^N denote the $N (\in \mathbb{N})$ dimensional Euclidean space, where \mathbb{N} denotes the set of nonnegative integers. The vector in \mathbb{R}^N with all components set to one is denoted by $\mathbf{1}_N$, and the vector with all components set to zero is denoted by $\mathbf{0}_N$. The identity matrix is denoted by \mathbf{I}_N , and the identity operator is denoted by I . The $N \times N$ matrix with all zero components

is denoted by \mathbf{O}_N . A k sparse vector $\mathbf{x} \in \mathbb{R}^N$ is a vector with at most k nonzero components. The i th element of a vector $\mathbf{x} \in \mathbb{R}^N$ is denoted by x_i , $i \in \{1, 2, \dots, N\}$. Similarly, the (i, j) component of a matrix $\mathbf{X} \in \mathbb{R}^{N \times M}$ is denoted by x_{ij} , where $i, j \in \mathbb{N}$ and $i \in \{1, 2, \dots, N\}$, $j \in \{1, 2, \dots, M\}$. Let $\text{tr}(\mathbf{A})$ denote the sum of the main diagonal elements of a square matrix. The ℓ_p norm ($p \in \{1, 2, \dots\}$) of vector $\mathbf{x} \in \mathbb{R}^N$ is defined by $\|\mathbf{x}\|_p := (\sum_{i=1}^N x_i^p)^{1/p}$. The inner product of $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$ is defined by $\langle \mathbf{x}, \mathbf{y} \rangle := \sum_{i=1}^N x_i y_i$. The Frobenius norm of a matrix $\mathbf{A} \in \mathbb{R}^{N \times M}$ is defined by $\|\mathbf{A}\|_F := (\sum_{i=1}^N \sum_{j=1}^M a_{ij}^2)^{1/2}$.

B. Convex analytic tools

A function $f : \mathbb{R}^N \rightarrow (-\infty, +\infty)$ is convex if $f(t\mathbf{x} + (1-t)\mathbf{y}) \leq tf(\mathbf{x}) + (1-t)f(\mathbf{y})$ for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$ and every $t \in (0, 1)$. A function is η -weakly convex if $f + \eta \frac{1}{2} \|\cdot\|_2^2$ is convex for some $\eta > 0$. Furthermore, a function is η -strongly convex when $f - \eta \frac{1}{2} \|\cdot\|_2^2$ is convex. We denote by $\Gamma_0(\mathbb{R}^N)$ the class of proper lower semicontinuous convex functions¹ from \mathbb{R}^N to $(-\infty, +\infty]$. The Fenchel conjugate function of $f \in \Gamma_0(\mathbb{R}^N)$ is defined by $f^*(\mathbf{x}) := \sup_{\mathbf{y} \in \mathbb{R}^N} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{y}))$ [22]. The Moreau envelope ${}^\gamma f$ of f of index $\gamma > 0$ is defined by ${}^\gamma f(\mathbf{x}) := \inf_{\mathbf{y} \in \mathbb{R}^N} (f(\mathbf{y}) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{y}\|_2^2)$ [22]. The proximity operator prox_f is defined by $\text{prox}_f(\mathbf{x}) := \arg \min_{\mathbf{y} \in \mathbb{R}^N} (f(\mathbf{y}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2)$ [22]. The smallest and largest eigenvalues of a symmetric matrix \mathbf{A} are denoted by $\lambda_{\min}(\mathbf{A})$ and $\lambda_{\max}(\mathbf{A})$.

C. Distributed sparse optimization

Let $\mathcal{G}(\mathcal{V}, \mathcal{E})$ be an undirected graph, where $\mathcal{V} := \{1, 2, \dots, m\}$ denotes a set of nodes (vertices), and \mathcal{E} denotes a set of edges. Here, $(i, j) \in \mathcal{E}$ means that there is an edge between $i \in \mathcal{V}$ and $j \in \mathcal{V}$. Each node $i \in \mathcal{V}$ is connected to $r_i = |\mathcal{N}_i|$ nodes, where $|\mathcal{N}_i|$ denotes the cardinality of the set \mathcal{N}_i of neighboring nodes of node i . We use the convention that $i \notin \mathcal{N}_i$. The graph is assumed connected; i.e., there exists a path that connects any pair of nodes in a single or multihop way. We consider a situation where each node i has the local sensing matrix $\mathbf{U}_i \in \mathbb{R}^{N \times l}$ and the local measurement vector $\mathbf{d}_i = \mathbf{U}_i^T \mathbf{w}^* + \mathbf{n}_i \in \mathbb{R}^l$, where $\mathbf{w}^* \in \mathbb{R}^N$ is the sparse unknown vector (which is common to every node), and $\mathbf{n}_i \in \mathbb{R}^l$, $i = 1, 2, \dots, m$, is the additive noise. The objective of distributed sparse optimization is to estimate the sparse vector \mathbf{w}^* over the entire network via local updates at each node and information exchanges among nodes.

III. PROPOSED APPROACH

We present the proposed problem formulation for distributed sparse optimization based on the consensus promoting penalty, and we then discuss the convexity of the cost function. We finally present a distributed optimization algorithm to solve the proposed formulation.

¹A convex function $f : \mathbb{R}^N \rightarrow (-\infty, +\infty]$ is proper if $f(\mathbf{x}) < +\infty$ for some $\mathbf{x} \in \mathbb{R}^N$. A convex function $f : \mathbb{R}^N \rightarrow (-\infty, +\infty]$ is lower semicontinuous if $\text{lev}_{\leq a} f := \{\mathbf{x} \in \mathbb{R}^N \mid f(\mathbf{x}) \leq a\}$ is a closed set for any $a \in \mathbb{R}$.

A. Formulation based on MC penalty

We start with the following penalized formulation:

$$(P_0) \min_{\mathbf{w} \in \mathbb{R}^N} \psi(\mathbf{w}) := \sum_{i=1}^m \underbrace{\left[\phi_i(\mathbf{w}) + \frac{\mu}{m} \pi(\mathbf{w}) \right]}_{=: \psi_i(\mathbf{w})}, \quad (1)$$

where we refer to $\psi(\mathbf{w})$ as the *vector global cost* and $\psi_i(\mathbf{w})$ as the *vector local cost* of problem (P₀) because its argument \mathbf{w} is a vector. Here, $\phi_i(\mathbf{w}) = \frac{1}{2} \|\mathbf{U}_i^T \mathbf{w} - \mathbf{d}_i\|_2^2$, $\mu > 0$, and

$$\pi(\mathbf{w}) := \|\mathbf{w}\|_1 - \min_{\mathbf{v} \in \mathbb{R}^N} \left(\|\mathbf{v}\|_1 + \frac{1}{2\gamma} \|\mathbf{w} - \mathbf{v}\|_2^2 \right) \quad (2)$$

is a weakly convex penalty called the minimax concave (MC) penalty [10], [11], where $\gamma \in (0, +\infty]$. The ℓ_1 norm is separable and proximable (in the sense that the proximity operator can be easily computed). To preserve the overall convexity of the entire cost, we assume that $\eta := \lambda_{\min}(\mathbf{U}\mathbf{U}^T) > 0$, so that $\phi(\mathbf{w}) := \sum_{i=1}^m \phi_i(\mathbf{w})$ is η -strongly convex for $\mathbf{U} := [\mathbf{U}_1 \ \mathbf{U}_2 \ \dots \ \mathbf{U}_m]$. In contrast, each ϕ_i is not strongly convex in typical scenarios.

By virtue of Moreau's decomposition [22] ${}^\gamma f + {}^{1/\gamma}(f^*) \circ \gamma^{-1}I = \frac{1}{2\gamma} \|\cdot\|_2^2$ of $f \in \Gamma_0(\mathbb{R}^N)$, the penalty can be rewritten as

$$\pi(\mathbf{w}) = \|\mathbf{w}\|_1 - \gamma^{-1} \left(\frac{1}{2} \|\mathbf{w}\|_2^2 - \gamma \frac{1}{\gamma} (\|\cdot\|_1^*)(\gamma^{-1}\mathbf{w}) \right). \quad (3)$$

Substituting (3) into (1) yields

$$\psi(\mathbf{w}) = \left(\sum_{i=1}^m \phi_i(\mathbf{w}) \right) - \frac{\mu}{2\gamma} \|\mathbf{w}\|_2^2 + \mu \|\mathbf{w}\|_1 + \mu \frac{1}{\gamma} (\|\cdot\|_1^*)(\gamma^{-1}\mathbf{w}). \quad (4)$$

Due to the η -strong convexity of $\phi(\mathbf{w})$ and the convexity of the third and fourth terms of (4), the convexity of the entire function ψ is ensured as long as $\eta - \mu/\gamma \geq 0$. However, the local function $\psi_i(\mathbf{w})$ cannot be convex unless $\phi_i(\mathbf{w})$ is strongly convex. Let $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_m]^T \in \mathbb{R}^{m \times N}$, where $\mathbf{w}_i \in \mathbb{R}^N$ is the estimate of $\mathbf{w}_{\text{opt}} := \arg \min_{\mathbf{w} \in \mathbb{R}^N} \psi(\mathbf{w})$ at node $i \in \mathcal{V}$. Note here that the vector global cost $\psi(\mathbf{w})$ possesses a unique minimizer due to the strong convexity of $\psi(\mathbf{w})$. Define the consensus subspace $\mathcal{C} := \{\mathbf{W} \in \mathbb{R}^{m \times N} \mid \mathbf{w}_1 = \mathbf{w}_2 = \dots = \mathbf{w}_m\}$. Problem (P₀) can then be formulated equivalently as the following constrained optimization problem:

$$(P_1) \min_{\mathbf{W} \in \mathcal{C}} \Psi(\mathbf{W}) := \sum_{i=1}^m \underbrace{\left[\phi_i(\mathbf{w}_i) + \frac{\mu}{m} \pi(\mathbf{w}_i) \right]}_{=: \psi_i(\mathbf{w}_i)}, \quad (5)$$

where we refer to $\Psi(\mathbf{W})$ as the *matrix global cost* and $\psi_i(\mathbf{w}_i)$ as the *matrix local cost* of problem (P₁) since the argument \mathbf{W} of Ψ is a matrix. Under consideration is the challenging case where $\Psi(\mathbf{W})$ is nonconvex. With weak convexity of $\frac{\mu}{m} \Pi(\mathbf{W}) := \sum_{i=1}^m \frac{\mu}{m} \pi(\mathbf{w}_i)$, the *matrix global cost* $\Psi(\mathbf{W})$ can be convex only if $\Phi(\mathbf{W}) := \sum_{i=1}^m \phi_i(\mathbf{w}_i)$ is strongly convex, which can only be true, due to its separability, if every term $\phi_i(\mathbf{w}_i)$ is strongly convex. This implies that the *matrix global cost* $\Psi(\mathbf{W})$ is convex only if every node has sufficient information that admits a unique solution. In the following subsection, we present a simple and effective method that

Table I: Types of costs for MC-CPP

	local (each node)	global (entire network)
vector \mathbf{w}	(i) $\psi_i(\mathbf{w})$ vector local cost	(ii) $\psi(\mathbf{w})$ vector global cost
matrix \mathbf{W}	(iii) $\theta_i(\mathbf{w}_i)$ matrix local cost	(iv) $\Theta(\mathbf{W})$ matrix global cost

allows to deal with the challenging, yet realistic, nonconvex case if $l < N$ (i.e., when the number l of local linear equations is smaller than the number N of local variables). In fact, $\Psi(\mathbf{W})$ can be convex only if $l \geq N$ is satisfied.²

B. Proposed formulation with consensus promoting penalty: convexity analysis

We introduce a consensus promoting penalty (CPP) term (which is inspired from the incidence matrix [23], [24]) $C(\mathbf{W}) := \sum_{i=1}^m \frac{\kappa_i}{2} \|\mathbf{w}_i - \bar{\mathbf{w}}_i\|_2^2$ to get the following minimax concave CPP (MC-CPP) formulation:

$$(P_2) \min_{\mathbf{W} \in \mathcal{C}} \Theta(\mathbf{W}) := \Psi(\mathbf{W}) + C(\mathbf{W}) \quad (6)$$

$$= \sum_{i=1}^m \underbrace{\left[\phi_i(\mathbf{w}_i) + \frac{\mu}{m} \pi(\mathbf{w}_i) + \frac{\kappa_i}{2} \|\mathbf{w}_i - \bar{\mathbf{w}}_i\|_2^2 \right]}_{\theta_i(\mathbf{w}_i)}, \quad (7)$$

where $\bar{\mathbf{w}}_i = \sum_{j \in \mathcal{N}_i} \alpha_{ij} \mathbf{w}_j$, $\kappa_i \geq 0, \forall i \in \{1, 2, \dots, m\}$, with $\alpha_{ij} = \alpha_{ji} \in (0, 1)$ such that $\sum_{j \in \mathcal{N}_i} \alpha_{ij} = 1$. Substituting (2) into (7), $\Theta(\mathbf{W})$ can be separated into smooth and nonsmooth parts as $\Theta(\mathbf{W}) = S(\mathbf{W}) + H(\mathbf{W})$, where

$$S(\mathbf{W}) = \underbrace{\sum_{i=1}^m \left(\phi_i(\mathbf{w}_i) + C_i(\mathbf{W}) \right)}_{=: F(\mathbf{W})} - \frac{\mu}{m\gamma} \sum_{i=1}^m \frac{1}{2} \|\mathbf{w}_i\|_2^2 + \frac{\mu}{m} \sum_{i=1}^m (\| \cdot \|_1^{1/\gamma}) (\gamma^{-1} \mathbf{w}_i), \quad (8)$$

and $H(\mathbf{W}) := \frac{\mu}{m} \sum_{i=1}^m \|\mathbf{w}_i\|_1$. Under certain parameter conditions, convexity of the matrix local cost can be ensured. However, this does not directly mean that the matrix global cost is convex. This is because CPP is regarded as a function of \mathbf{w}_i in the matrix local cost, while it is a function of \mathbf{W} in the matrix global cost. We summarize the four types of cost function in Table I. The vector global cost of the MC-CPP is $\psi(\mathbf{w})$ because $C(\mathbf{W}) = 0$ for any $\mathbf{W} \in \mathcal{C}$. Likewise, the vector local cost is $\psi_i(\mathbf{w})$ because $\theta_i(\mathbf{w}) = \psi_i(\mathbf{w})$ for any $\mathbf{W} \in \mathcal{C}$. Since $\mathbf{w}_i = \mathbf{W}^\top \mathbf{e}_i \in \mathbb{R}^N$, we rewrite CPP as $C_i(\mathbf{W}) = \frac{\kappa_i}{2} \|\mathbf{w}_i - \bar{\mathbf{w}}_i\|_2^2 = \frac{\kappa_i}{2} \|\mathbf{W}^\top \mathbf{v}_i\|_2^2$, where $\mathbf{v}_i := \mathbf{e}_i - \sum_{j \in \mathcal{N}_i} \alpha_{ij} \mathbf{e}_j \in \mathbb{R}^m$. Let us define $\mathbf{W} =: [\hat{\mathbf{w}}_1 \ \hat{\mathbf{w}}_2 \ \dots \ \hat{\mathbf{w}}_m]$,

²In environmental modeling scenarios, for instance, weather sensors may take few measurements and the environment changes slowly compared with the time to apply many iterations of the algorithm.

Table II: Convexity condition for each cost with MC-CPP

	local	global
vector \mathbf{w}	(i) $\eta_i > \frac{\mu}{m\gamma}$ (but $\eta_i = 0$ typically)	(ii) $\mu/\gamma < \eta$
matrix \mathbf{W}	(iii) $\kappa \geq \frac{\mu}{m\gamma}$	(iv) To be investigated (see Figure 2)

and $\mathbf{V}_\kappa := [\sqrt{\kappa_1} \mathbf{v}_1 \ \sqrt{\kappa_2} \mathbf{v}_2 \ \dots \ \sqrt{\kappa_m} \mathbf{v}_m] \in \mathbb{R}^{m \times m}$. The CPP term can then be expressed as

$$C(\mathbf{W}) := \sum_{i=1}^m C_i(\mathbf{W}) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^N \kappa_i \hat{\mathbf{w}}_j^\top \mathbf{v}_i \mathbf{v}_i^\top \hat{\mathbf{w}}_j$$

$$= \frac{1}{2} \sum_{j=1}^N \hat{\mathbf{w}}_j^\top \mathbf{V}_\kappa \mathbf{V}_\kappa^\top \hat{\mathbf{w}}_j. \quad (9)$$

Strong convexity of the consensus promoting penalty can thus be analyzed by inspecting strict positivity of the eigenvalues of $\mathbf{V}_\kappa \mathbf{V}_\kappa^\top$. We now highlight the following equivalence:

$$C(\mathbf{W}) = 0 \Leftrightarrow \mathbf{W} \in \mathcal{C}. \quad (10)$$

Here, \Leftrightarrow is clear. To verify the converse \Rightarrow , we express \mathbf{V}_κ as $\mathbf{V}_\kappa = \mathbf{V} \mathbf{\Lambda}_\kappa$ with $\mathbf{V} := [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_m] \in \mathbb{R}^{m \times m}$, and $\mathbf{\Lambda}_\kappa := \text{diag}(\sqrt{\kappa_1}, \sqrt{\kappa_2}, \dots, \sqrt{\kappa_m})$, where $\text{diag}(\cdot)$ is the diagonal matrix with its diagonal entries given by the arguments. It is then sufficient, from the nonsingularity of the matrix $\mathbf{\Lambda}_\kappa$, to show that $(\text{Ker} \mathbf{V}_\kappa^\top =) \text{Ker} \mathbf{V}^\top = \text{span}\{\mathbf{1}\}$. Clearly, the matrix $\mathbf{Y}^\top := \mathbf{I}_m - \mathbf{V}^\top$ shares the eigenspace with \mathbf{V}^\top , and the eigenvalue one of \mathbf{Y}^\top and the eigenvalue zero of \mathbf{V}^\top both correspond to the same eigenvector $\mathbf{1}_m$. Now, \mathbf{Y}^\top is a nonnegative matrix that is irreducible³ as the graph is assumed connected. The classical Perron Frobenius theory thus tells us that the eigenvalue one corresponding to the positive eigenvector $\mathbf{1}_m$ is simple (has algebraic multiplicity one), which implies the equality of the assertion. By (10), $C(\mathbf{W})$ does not change the entire cost on the consensus subspace \mathcal{C} , meaning that the global minimizer (the optimal solution) is preserved. More importantly, it also implies that the function $C(\mathbf{W})$ is nonzero, and it is actually strongly convex on the orthogonal complement \mathcal{C}^\perp of \mathcal{C} . We summarize the important points below. (See [1] for proofs of the lemmas and proposition given in the following.)

- Lemma 1.** (a) Each matrix local cost $\theta_i(\mathbf{w}_i)$ is convex if $\kappa_i \geq \frac{\mu}{m\gamma}$.
(b) Each vector local cost $\psi_i(\mathbf{w})$ is convex iff $(\lambda_{\min}(\mathbf{U}_i \mathbf{U}_i^\top) =: \eta_i) \geq \frac{\mu}{m\gamma}$.
(c) The vector global cost $\psi(\mathbf{w})$ is convex iff $(\lambda_{\min}(\mathbf{U} \mathbf{U}^\top) =: \eta) \geq \frac{\mu}{\gamma}$.
(d) $\min_{\mathbf{W} \in \mathcal{C}} \Theta(\mathbf{W}) = \min_{\mathbf{W} \in \mathcal{C}} \Psi(\mathbf{W})$.
(e) $\arg \min_{\mathbf{W} \in \mathcal{C}} \Theta(\mathbf{W}) = \arg \min_{\mathbf{W} \in \mathcal{C}} \Psi(\mathbf{W})$.

The convexity conditions are summarized in Table II. Note here that $\eta_i = 0$ (i.e., ϕ_i is not strongly convex) in typical situations, because each local node is supposed to have no sufficient amount of measurements to identify the solution

³A square matrix is said to be irreducible if the matrix cannot be reorganized into a block upper-triangular form by simultaneous row/column permutations.

so the nodes need to collaborate with each other (cf. Section III-A).

Lemma 2. *Regarding the costs (i) – (iv) presented in Table I, the following implications hold.*

- 1) Convexity of (i) implies convexity of (ii), (iii), (iv).
- 2) Convexity of (iv) implies convexity of (ii).

Proposition 1. *Let range $\mathbf{U} = \mathbb{R}^N$ so that the function $\mathbf{w} \mapsto \sum_{i=1}^m \phi_i(\mathbf{w})$ is η -strongly convex for $\eta := \lambda_{\min}(\mathbf{U}\mathbf{U}^\top)$ (see Section III-A). Assume that the graph is connected. Then $\Theta(\mathbf{W})$ is convex only if $\mu \leq \eta\gamma$.*

Most existing algorithms for distributed convex optimization assume convexity of the *vector local cost*, which is not apparent with the use of CPP. In some algorithms such as the PG-EXTRA, however, the convexity of the vector local cost itself is not used explicitly to prove the convergence of the algorithm, and it is only used to ensure convexity of the matrix global cost, which is used directly to prove the convergence. Therefore, convexity of the matrix global cost will be studied by simulations in Section IV-B.

C. Distributed proximal debiasing-gradient (DPD) method

Let $\mathbf{W}_k \in \mathbb{R}^{m \times N}$ denote the variable matrix at time $k \in \mathbb{N}$, and $\mathbf{M}, \tilde{\mathbf{M}} \in \mathbb{R}^{m \times m}$ are the mixing matrices.⁴ The gradient of $S(\mathbf{W})$ is denoted by $\nabla S(\mathbf{W}) = [\nabla s_1(\mathbf{w}_1) \nabla s_2(\mathbf{w}_2) \dots \nabla s_m(\mathbf{w}_m)]^\top$. By using the property $\nabla^\gamma f = \gamma^{-1}(I - \text{prox}_{\gamma f})$ for $f \in \Gamma_0(\mathbb{R}^N)$, $\gamma > 0$ [22, Ch. 14] in (8), we have $\nabla s_i(\mathbf{w}) = \nabla \phi_i(\mathbf{w}) - \frac{\mu}{\gamma m}(\mathbf{w} - \text{prox}_{\gamma h}(\mathbf{w})) + \kappa_i(\mathbf{w} - \bar{\mathbf{w}}_i)$. By the triangle inequality and the (firm) nonexpansivity of $\text{prox}_{\gamma h}$, the following inequality holds:

$$\begin{aligned} & \|\nabla s_i(\mathbf{x}) - \nabla s_i(\mathbf{y})\|_2 \\ & \leq \|\nabla \phi_i(\mathbf{x}) - \nabla \phi_i(\mathbf{y})\|_2 + \left\| \left(\kappa_i - \frac{\mu}{\gamma m} \right) (\mathbf{x} - \mathbf{y}) \right\|_2 \\ & + \frac{\mu}{\gamma m} \left\| \text{prox}_{\gamma \|\cdot\|_1}(\mathbf{x}) - \text{prox}_{\gamma \|\cdot\|_1}(\mathbf{y}) \right\|_2 \leq (L_i + \kappa_i) \|\mathbf{x} - \mathbf{y}\|_2. \end{aligned}$$

Here, $L_i := \lambda_{\max}^{1/2}(\mathbf{U}_i \mathbf{U}_i^\top)$ is the Lipschitz constant of $\nabla \phi_i$; i.e., $\|\nabla \phi_i(\mathbf{x}) - \nabla \phi_i(\mathbf{y})\|_2 \leq L_i \|\mathbf{x} - \mathbf{y}\|_2$ for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$.

An application of PG-EXTRA [7] to the proposed formulation (P₂) yields Algorithm 1. Note that other algorithms can be applied to (P₂) that assume convex, proximable penalties. We refer to this particular method utilizing the two proximity operators as the distributed proximal debiasing-gradient method (DPD). Hereafter, we assume that all nodes have the same number of neighbors i.e., $r_1 = r_2 = \dots = r_m =: r$ for simplicity.

IV. NUMERICAL EXAMPLES

We study the performance of the proposed algorithm for $\phi_i(\mathbf{w}_i) = \frac{1}{2} \|\mathbf{u}_i^\top \mathbf{w}_i - d_i\|_2^2$ ($l = 1$), where $\mathbf{u}_i \in \mathbb{R}^N$ is the input vector, distributed according to $\mathbf{u}_i \sim \mathcal{N}(\mathbf{0}_N, \mathbf{I}_N)$ $i = 1, 2, \dots, m$, and $d_i = \mathbf{u}_i^\top \mathbf{w}^* + n_i$ is the output (see Section 2.2). The normalized error is defined as $\text{Error}(\mathbf{x}, \mathbf{y}) =$

⁴The mixing matrices \mathbf{M} and $\tilde{\mathbf{M}}$ are subject to the following conditions [7]: (i) If $i \neq j$ and $(i, j) \notin \mathcal{E}$, then $M_{ij} = \tilde{M}_{ij} = 0$, (ii) $\mathbf{M} = \mathbf{M}^\top$, $\tilde{\mathbf{M}} = \tilde{\mathbf{M}}^\top$, (iii) $\text{null}(\mathbf{M} - \tilde{\mathbf{M}}) = \text{span}\{\mathbf{1}_m\}$, $\text{null}(\mathbf{I}_m - \tilde{\mathbf{M}}) \supseteq \text{span}\{\mathbf{1}_m\}$, and (iv) $\tilde{\mathbf{M}} \succ 0$ and $\frac{\mathbf{I}_m + \tilde{\mathbf{M}}}{2} \succeq \tilde{\mathbf{M}} \succeq \mathbf{M}$.

Algorithm 1 DPD-CPP derived from PG-EXTRA

Require: $\kappa_{i_z} > 0$, $i = 1, 2, \dots, m$, $0 < \beta \leq 2\lambda_{\min}(\tilde{\mathbf{M}})/\max\{L_i + \kappa_i\}_{i=1}^m$, $\mu > 0$ and $\gamma > 0$ such that $\mu/\gamma < \min\{\lambda_2(\mathbf{V}_\kappa \mathbf{V}_\kappa^\top), \eta\}$
 \mathbf{W}_0 is an arbitrary point

- 1: $\nabla S(\mathbf{W}_0) = \nabla \Phi(\mathbf{W}_0) - \frac{\mu}{m\gamma} \mathbf{W}_0 + \text{diag}(\kappa_1, \kappa_2, \dots, \kappa_m)(\mathbf{W}_0 - \bar{\mathbf{W}}_0) + \frac{\mu}{m\gamma} \text{prox}_{\gamma \|\cdot\|_1}(\mathbf{W}_0)$
- 2: $\mathbf{W}_{\frac{1}{2}} = \mathbf{M}\mathbf{W}_0 - \beta \nabla S(\mathbf{W}_0)$
- 3: $\mathbf{W}_1 = \text{prox}_{\frac{\mu\beta}{m} \|\cdot\|_1}(\mathbf{W}_0)$
- 4: **for** $k = 0, 1, 2, \dots$ **do**
- 5: $\nabla S(\mathbf{W}_{k+1}) = \nabla \Phi(\mathbf{W}_{k+1}) - \frac{\mu}{m\gamma} \mathbf{W}_{k+1} + \text{diag}(\kappa_1, \kappa_2, \dots, \kappa_m)(\mathbf{W}_{k+1} - \bar{\mathbf{W}}_{k+1}) + \frac{\mu}{m\gamma} \text{prox}_{\gamma \|\cdot\|_1}(\mathbf{W}_{k+1})$
- 6: $\mathbf{W}_{k+1+\frac{1}{2}} = \mathbf{M}\mathbf{W}_{k+1} + \mathbf{W}_{k+\frac{1}{2}} - \tilde{\mathbf{M}}\mathbf{W}_k - \beta [\nabla S(\mathbf{W}_{k+1}) - \nabla S(\mathbf{W}_k)]$
- 7: $\mathbf{W}_{k+2} = \text{prox}_{\frac{\mu\beta}{m} \|\cdot\|_1}(\mathbf{W}_{k+1+\frac{1}{2}})$
- 8: **end for**

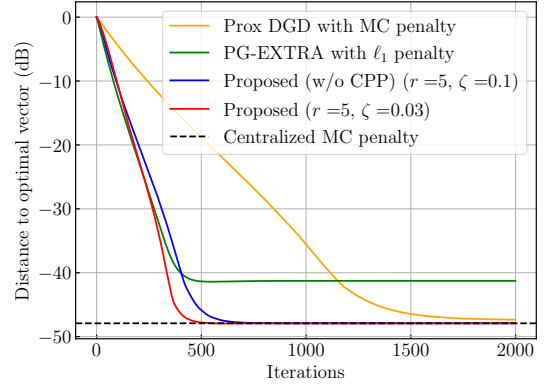


Figure 1: Learning curves in system mismatch of the average estimates.

$10 \log_{10} \left(\frac{\|\mathbf{x} - \mathbf{y}\|_2^2}{\|\mathbf{x}\|_2^2} \right)$. Furthermore, the components of the mix-

ing matrix \mathbf{M} are set to $m_{ij} = \begin{cases} \frac{\zeta}{r+1} & \text{if } j \in \mathcal{N}_i \\ 1 - \frac{\zeta r}{r+1} & \text{if } i = j \\ 0 & \text{otherwise,} \end{cases}$

$\zeta \in (0, \frac{r+1}{r})$ indicates the strength of the influence of the neighboring estimates, and $\tilde{\mathbf{M}} = \frac{\mathbf{I}_m + \mathbf{M}}{2}$.

A. Learning speed of proposed approaches

The dimension of \mathbf{w}^* is $N = 10$, the network consists of $m = 100$ nodes with $r = 5$ neighboring nodes, and the signal to noise ratio (SNR) is 30 dB. The estimandum \mathbf{w}^* is 30% sparse. For CPP, we let $\mathbf{V}^\top = \frac{r+1}{r}(\mathbf{I}_m - \mathbf{M})$ with $\zeta = 1$. The parameters of the proposed algorithm are set to $\mu = 9.0 \times 10^{-3}$, $\gamma = 0.15$, $\beta = 0.022$, $(\zeta, \kappa) = (0.15, 0), (0.4, 10)$. The parameters of PG-EXTRA with the ℓ_1 penalty (DPD with $\gamma = +\infty, \kappa = 0$) are set to $\mu = 1.9 \times 10^{-3}$, $\beta = 0.011$, and $\zeta = 1$. Prox-DGD [17] is also tested for reference, where the step size of Prox-DGD is set to 0.014, and the parameters for the MC penalty are set to the same value as DPD.

Figure 1 plots $\text{Error}(\mathbf{w}^*, \frac{1}{m} \mathbf{W}^\top \mathbf{1}_m)$, the system mismatch of an average estimate over nodes, where the dashed line

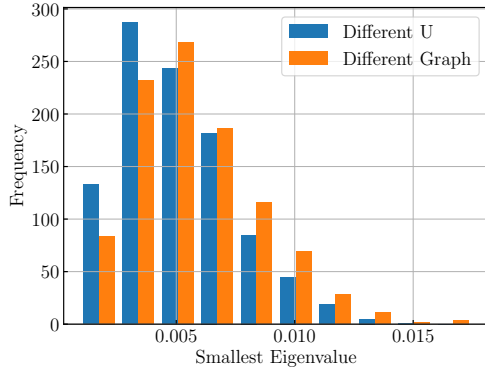


Figure 2: A histogram of the smallest eigenvalue of Hessian of $\Theta(\mathbf{W})$ with different \mathbf{U} and different graphs.

at the bottom is a benchmark, showing the performance of a centralized approach to problem (P_0) . One can see that the proposed MC-based approaches significantly outperform the ℓ_1 -based approach, while the Prox-DGD method suffers from slow convergence. One can also see that the proposed method achieves the same error as the centralized solution (the benchmark) even without CPP, which leaves the possibility that PG-EXTRA applied to the minimization problem of $\Psi(\mathbf{W})$ would still enjoy provable convergence.

B. Convexity with CPP

We study the convexity of the cost function with the CPP term, which is designed for the cost to become convex when $\mathbf{W} \in \mathcal{C}^\perp$. However, the assurance of the convexity over the whole space $\mathbb{R}^{N \times m}$ still remains to be proven. Specifically, we inspect strong convexity of the quadratic function $\Phi(\mathbf{W}) + C(\mathbf{W})$ for $\mu = 0$, or equivalently strict positivity of the smallest eigenvalue of its Hessian matrix. The parameters are set to $N = 10$, $m = 50$, and $r = 2$.

Figure 2 plots the frequency of the smallest eigenvalue of Θ over 1000 trials. The blue bar in the figure considers the case when \mathbf{U} is changed using `numpy.random.randn`, while the orange bar considers the case when the weights of the graph are fixed to $r = \frac{1}{d}$ with the links of underlying graph selected randomly. In both cases, it is clear that the smallest eigenvalue remains positive. This means that, in this scenario, CPP would serve to convexify the matrix global cost, which may suggest that weakly convex penalties can be used in distributed settings with convergence guarantees to a global minimizer.

V. CONCLUDING REMARKS

We introduce a consensus promoting penalty that has the object of convexifying a cost function containing the weakly convex MC penalty. The convexity of the cost is desirable because the formulation can yield accurate sparse estimates with convergence guarantees to a global minimizer. Numerical examples showed that the CPP convexified the entire matrix global cost, and it also accelerated the speed of convergence. The use of the MC penalty in online scenarios has been studied in [25], and an extension of the present work to online settings will be an interesting future work.

REFERENCES

- [1] K. Komuro, M. Yukawa, and R. L. G. Cavalcante. Distributed sparse optimization with weakly convex regularizer: Consensus promoting and approximate moreau enhanced penalties towards global optimality. *IEEE Transactions on Signal and Information Processing over Networks*, 2022. accepted for publication.
- [2] D. Bertsekas and J. Tsitsiklis. *Parallel and Distributed Computation: Numerical Methods*. Prentice-Hall, Inc., USA, 1989.
- [3] T. Chang, M. Hong, H. Wai, X. Zhang, and S. Lu. Distributed learning in the nonconvex world: From batch data to streaming and beyond. *IEEE Signal Processing Magazine*, 37(3):26–38, 2020.
- [4] I. Notarnicola, Y. Sun, G. Scutari, and G. Notarstefano. Distributed big-data optimization via blockwise gradient tracking. *IEEE Transactions on Automatic Control*, 66(5):2045–2060, 2021.
- [5] Q. Ling and Z. Tian. Decentralized sparse signal recovery for compressive sleeping wireless sensor networks. *IEEE Transactions on Signal Processing*, 58(7):3816–3827, 2010.
- [6] A. I. Chen and A. Ozdaglar. A fast distributed proximal-gradient method. In *2012 50th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 601–608, 2012.
- [7] W. Shi, Q. Ling, G. Wu, and W. Yin. A proximal gradient algorithm for decentralized composite optimization. *IEEE Transactions on Signal Processing*, 63(22):6013–6023, 2015.
- [8] J. Fan and R. Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *J. American Statistical Association*, 96(456):1348–1360, Dec. 2001.
- [9] M. Yukawa and S. Amari. ℓ_p -regularized least squares ($0 < p < 1$) and critical path. *IEEE Trans. Information Theory*, 62(1):488–502, Jan. 2016.
- [10] I. Selesnick. Sparse regularization via convex analysis. *IEEE Transactions on Signal Processing*, 65(17):4481–4494, 2017.
- [11] C. Zhang. Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics*, 38(2):894 – 942, 2010.
- [12] J. Abe, M. Yamagishi, and I. Yamada. Linearly involved generalized moreau enhanced models and their proximal splitting algorithm under overall convexity condition. *Inverse Problems*, 36(3):1–36, Nov. 2019.
- [13] M. Yukawa, H. Kaneko, K. Suzuki, and I. Yamada. Linearly-involved Moreau-enhanced-over-subspace model: debiased sparse modeling and stable outlier-robust regression. 2021. [Online]. Available: <https://arxiv.org/abs/2201.03235>.
- [14] P. Lorenzo and G. Scutari. Distributed nonconvex optimization over time-varying networks. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4124–4128, 2016.
- [15] D. Hajinezhad and M. Hong. Nonconvex alternating direction method of multipliers for distributed sparse principal component analysis. In *2015 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pages 255–259, 2015.
- [16] H. Sun and M. Hong. Distributed non-convex first-order optimization and information processing: Lower complexity bounds and rate optimal algorithms. In *2018 52nd Asilomar Conference on Signals, Systems, and Computers*, pages 38–42, 2018.
- [17] J. Zeng and W. Yin. On nonconvex decentralized gradient descent. *IEEE Transactions on Signal Processing*, 66(11):2834–2848, 2018.
- [18] P. D. Lorenzo and G. Scutari. Next: In-network nonconvex optimization. *IEEE Transactions on Signal and Information Processing over Networks*, 2(2):120–136, 2016.
- [19] A. Daneshmand, G. Scutari, and V. Kungurtsev. Second-order guarantees of distributed gradient algorithms. *SIAM Journal on Optimization*, 30(4):3029–3068, 2020.
- [20] M. Hong, M. Razaviyayn, and J. Lee. Gradient primal-dual algorithm converges to second-order stationary solution for nonconvex distributed optimization over networks. In *Proc. ICLM*, volume 80, pages 2009–2018, 2018.
- [21] K. Komuro, M. Yukawa, and R. L. G. Cavalcante. Distributed sparse optimization with minimax concave regularization. in *Proc. IEEE Statistical Signal Processing Workshop*, 2021.
- [22] H. Bauschke and P. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Space*. Springer, 2nd edition, 2017.
- [23] D. Hajinezhad M. Hong and M. Zhao. Prox-pda: The proximal primal-dual algorithm for fast distributed nonconvex optimization and learning over networks. in *Proc. of the 34th International Conference on Machine Learning*, pages 1529–1538, 2017.
- [24] L. Jian, Y. Zhao, J. Hu, and P. Li. Distributed inexact consensus-based admm method for multi-agent unconstrained optimization problem. *IEEE Access*, 7:79311–79319, 2019.
- [25] H. Kaneko and M. Yukawa. Normalized least-mean-square algorithms with minimax concave penalty. In *Proc. IEEE ICASSP*, pages 5445–5449, 2020.