# A UNIFIED DECENTRALIZED NONCONVEX ALGORITHM UNDER KURDYKA-ŁOJASIEWICZ PROPERTY

HAO WU, LIPING WANG, AND HONGCHAO ZHANG

ABSTRACT. In this paper, we study the decentralized optimization problem of minimizing a finite sum of continuously differentiable and possibly nonconvex functions over a fixed-connected undirected network. We propose a unified decentralized nonconvex algorithmic framework that subsumes existing state-of-the-art gradient tracking algorithms and particularly several quasi-Newton algorithms. We present a general analytical framework for the convergence of our unified algorithm under both nonconvex and the Kurdyka-Łojasiewicz condition settings. We also propose some quasi-Newton variants that fit into our framework, where economical implementation strategies are derived for ensuring bounded eigenvalues of Hessian inverse approximations. Our numerical results show that these newly developed algorithms are very efficient compared with other state-of-the-art algorithms for solving decentralized nonconvex smooth optimization.

## 1. INTRODUCTION

In the era of data explosion and connected intelligence, decentralized optimization has emerged as a fundamental computational paradigm for large-scale and privacy-aware systems. This paper focuses on solving optimization problems over multi-node networks where no central server exists-a setting that naturally arises in modern applications including but not limited to decentralized resources control [17], wireless networks [26], decentralized machine learning [64], power systems [21], federated learning [41].

We consider a network of $n$ nodes, interconnected via an undirected and connected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \{1, \ldots, n\}$ is the set of nodes, and $\mathcal{E}$ is the collection of unordered edges. Each node $i \in \mathcal{V}$ possesses a private local objective function $f_i : \mathbb{R}^p \to \mathbb{R}$ and the collective goal is to solve the consensus optimization problem:

$$(1.1) \qquad \min_{\mathbf{z} \in \mathbb{R}^p} F(\mathbf{z}) := \frac{1}{n} \sum_{i=1}^{n} f_i(\mathbf{z}),$$

where each $f_i$ is continuously differentiable, possibly nonconvex, and known only to node $i$. Unlike the distributed setup with a server, the fully decentralized architecture requires all nodes to collaboratively reach an optimal solution using only local computation and peer-to-peer communication.

Because of wide and important practical applications, decentralized optimization has been extensively studied, where gradient-based first-order methods have

attracted much attention due to their simple implementation and low computation cost at each iteration. Existing well-developed works include Decentralized Gradient Descent (DGD) [38, 58, 60], gradient tracking methods [57, 40, 36, 37] and their variants [25, 18, 63, 48], exact diffusion methods [44, 59, 2, 53], momentum methods [55, 50, 20], and primal-dual methods [45, 23, 33], etc.

The convergence rates of various gradient-based methods are highly sensitive to the condition number of the objective function. In practical applications, the considered objective could be badly scaled or ill-conditioned. The Newton type and quasi-Newton techniques are efficient ways to improve the convergence speed of algorithms. A large number of studies focus on leveraging Hessian matrices to design second-order decentralized algorithms [34, 35, 61, 24, 30]. Numerical experiments have demonstrated that these algorithms exhibit significantly faster convergence compared to gradient-based algorithms. However, despite their empirical success, these methods are theoretically limited to a no-faster-than-linear convergence rate and incur significant computation costs due to the calculation of Hessian matrices or their inverses. Quasi-Newton techniques are more widely used in decentralized optimization to promote quasi-Newton algorithm development [16, 29, 62, 46, 52], since the curvature information of the Hessian of the objective function can be easily captured with a low computation cost.

It should be noted that the existing decentralized quasi-Newton methods [16, 29, 62, 46, 52] are proposed for solving convex or strongly convex minimization problems. *It is unknown whether some of these methods can be applied to nonconvex optimization with convergence guarantee.* In centralized nonconvex settings, the convergence of quasi-Newton methods has been extensively studied with various correction techniques, which tackle non-convexity and ill-conditioning to ensure robust performance [27, 31, 9]. Remarkably, we find in this decentralized quasi-Newton method under the strongly convex setting [62] that the regularization and damping techinques are able to guarantee positive definiteness of the quasi-Newton update even though the problem is nonconvex. This observation leads us to believe a certain class of decentralized quasi-Newton methods could have the provable convergence for nonconvex optimization.

We notice that many works [22, 3, 49, 56, 8] have attempted to unify some of various gradient-based algorithms in the convex or strongly convex settings while few works [4, 15] provide a unified analysis of several well-known decentralized methods for nonconvex case. Notably, these framework [4, 15] canot cover any quasi-Newton method. In this paper, we unify and generalize a class of decentralized algorithms under the nonconvex settings. All quasi-Newton's variants in this framework have global convergence (stationarity of every limit point of the iterates) under mild assumptions and proper choice of step sizes.

In recent years, the Kurdyka-Łojasiewicz (KŁ) property, initially introduced by Łojasiewicz [32], has been successfully employed to strengthen the convergence analysis of many optimization algorithms for nonconvex and nonsmooth optimization problems in the centralized setting [7, 39]. However, extending such analyses to the decentralized setting poses significant challenges, primarily due to the non-monotone decrease of the objective function sequence, which complicates the full convergence analysis of iterates for nearly all decentralized algorithms. Several preliminary works have explored the convergence of decentralized algorithms under the KŁ property [60, 13, 10, 51]. Most of these studies [60, 13, 51] rely on an

abstract convergence framework that assumes the existence of a bounded sequence generated by a descent algorithm, along with a potential function that also satisfies the KŁ property. However, the KŁ property imposed on the potential function cannot generally transfer to the objective function, vice versa. Although Chen et al. [10] effectively leveraged the KŁ property directly on the objective function, they postulate a strong connectivity assumption on the underlying network. In this paper, we proceed to utilize the KŁ property to eastablish the refined whole sequence convergence for all algorithms subsumed into our proposed framework.

Our main contributions in this paper can be summarized as follows.

1. We introduce a Unified Decentralized Nonconvex Algorithm (UDNA) which is a novel framework that systematically incorporates various gradient tracking methods [57, 40, 36, 37, 14, 43] and several well-developed quasi-Newton methods [25, 18, 46, 62]. So, we can investigate and summarize existing decentralized algorithms, and design new efficient algorithms by choosing from different combinations of communication and second-order approximation strategies, all within a cohesive analytical structure.

2. We first establish a robust subsequence convergence theory for the UDNA family under standard assumptions (Lipschitz gradients and lower-bounded objectives), provided the approximate matrices remain positive definite with bounded eigenvalues. More significantly, we overcome the key challenge of nonmonotonicity in decentralized optimization by successfully incorporating the KŁ property of the objective function itself. This analysis not only guarantees convergence to stationary points but also, for the first time in the decentralized literature, establishes the full sequence convergence of iterates generated by a quasi-Newton based framework, a result previously unavailable for such a broad class of methods.

3. We instantiate the UDNA framework with several innovative and computationally efficient quasi-Newton variants. A key technical contribution is the design of novel Hessian inverse approximation techniques that inherently ensure bounded eigenvalues, making them directly applicable within our framework. These newly proposed methods are so economical that they require only vector-vector products and minimize per-iteration computation and memory costs. Our numerical experiments on nonconvex optimization problems show newly proposed methods are very effective compared with several state-of-the-art decentralized methods.

The paper is organized as follows. In Section 2, we first propose the UDNA framework, show its subsequence convergence, and provide the enhanced convergence analysis under the KŁ property. The approaches of generating Hessian inverse approximation are also presented in this section. Numerical experiments of comparing our new methods with other well-established methods for solving decentralized nonconvex optimization are presented in Section 3. We finally draw some conclusions in Section 4.

1.1. **Notation.** We use uppercase boldface letters, e.g. $\mathbf{W}$, for matrices and lowercase boldface letters, e.g. $\mathbf{w}$, for vectors. For any vectors $\mathbf{v}_i \in \mathbb{R}^p$, $i = 1, \ldots, n$, we define $\bar{\mathbf{v}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{v}_i$ and $\mathbf{v} = [\mathbf{v}_1; \mathbf{v}_2; \ldots; \mathbf{v}_n] \in \mathbb{R}^{np}$. Given an undirected network $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, let $\mathbf{x}_i$ denote the local copy of the global variable $\mathbf{z}$ at node $i$ and $\mathcal{N}_i$ denote the set consisting of the neighbors of node $i$ (for convenience, we

treat node $i$ itself as one of its neighbors). We define $f(\mathbf{x}) = \sum_{i=1}^{n} f_i(\mathbf{x}_i)$ and use $\mathbf{g}^t$, $\mathbf{g}_i^t$ to stand for $\nabla f(\mathbf{x}^t)$, $\nabla f_i(\mathbf{x}_i^t)$ respectively, where, for clarification, the gradient of $f(\mathbf{x})$ is defined as $\nabla f(\mathbf{x}) = [\nabla f_1(\mathbf{x}_1); \nabla f_2(\mathbf{x}_2); \ldots; \nabla f_n(\mathbf{x}_n)] \in \mathbb{R}^{np}$. In addition, we define $\overline{\nabla} f(\mathbf{x}^t) = \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(\mathbf{x}_i^t) \in \mathbb{R}^p$. We say that $\mathbf{x}$ is consensual or gets consensus if $\mathbf{x}_1 = \mathbf{x}_2 = \ldots = \mathbf{x}_n$. $\mathbf{I}_p$ denotes the $p \times p$ identity matrix and $\mathbf{I}$ denote $\mathbf{I}_{np}$ for simplicity. Kronecker Product is denoted as $\otimes$. Given a vector $\mathbf{v}$ and a symmetric matrix $\mathbf{N}$, $\mathrm{span}(\mathbf{v})$ stands for the linear subspace spanned by $\mathbf{v}$; $\mathrm{Null}(\mathbf{N})$ and $\mathbf{N}^{\mathsf{T}}$ denote the null space and transpose of $\mathbf{N}$, respectively; $\lambda_{\min}(\mathbf{N})$, $\lambda_{\max}(\mathbf{N})$, and $\rho(\mathbf{N})$ denote the smallest eigenvalue, the largest eigenvalue, and the spectral radius of $\mathbf{N}$, respectively; For symmetric matrices $\mathbf{N}_1$ and $\mathbf{N}_2$ with same dimension, $\mathbf{N}_1 \succeq \mathbf{N}_2$ means $\mathbf{N}_1 - \mathbf{N}_2$ is positive semidefinite, while $\mathbf{N}_1 \geq \mathbf{N}_2$ means $\mathbf{N}_1 - \mathbf{N}_2$ is component-wise nonnegative. We denote $\log_{10}(\cdot)$ by $\log(\cdot)$ and define $\mathbf{M} = \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^{\mathsf{T}} \otimes \mathbf{I}_p$ where $\mathbf{1}_n \in \mathbb{R}^n$ denotes the vector with all components ones.

## 2. ALGORITHM DEVELOPMENT, CONVERGENCE ANALYSIS, AND PRACTICAL IMPLEMENTATION

The following are several necessary assumptions for the objective function.

**Assumption 1.** *Every local objective function $f_i$ is proper (i.e., not everywhere infinite) and coercive (a function $h$ is called coercive if $\|\mathbf{z}\| \to +\infty$ implies $h(\mathbf{z}) \to +\infty$).*

**Assumption 2.** *Each local gradient $\nabla f_i$ is Lipschitz continuous with constant $L > 0$, i.e.,*

$$(2.1) \qquad \|\nabla f_i(\mathbf{z}) - \nabla f_i(\tilde{\mathbf{z}})\| \leq L \|\mathbf{z} - \tilde{\mathbf{z}}\|,$$

$\forall \mathbf{z}, \tilde{\mathbf{z}} \in \mathbb{R}^p, i \in \mathcal{V}$.

Each $f_i$ is lower bounded following Assumption 1. In decentralized optimization it is convenient to parameterize communication by a mixing matrix $\tilde{\mathbf{W}} = [\tilde{W}_{ij}] \in \mathbb{R}^{n \times n}$, which is defined as follows.

**Definition 2.1.** (Mixing matrix $\tilde{\mathbf{W}}$ for given network $\mathcal{G} = (\mathcal{V}, \mathcal{E})$)

1. $\tilde{\mathbf{W}}$ is nonnegative, where each component $\tilde{W}_{ij}$ characterizes the active link $(i, j)$, i.e., $\tilde{W}_{ij} > 0$ if $j \in \mathcal{N}_i$; $\tilde{W}_{ij} = 0$, otherwise.
2. $\tilde{\mathbf{W}}$ is symmetric and doubly stochastic, i.e., $\tilde{\mathbf{W}} = \tilde{\mathbf{W}}^{\mathsf{T}}$ and $\tilde{\mathbf{W}} \mathbf{1}_n = \mathbf{1}_n$.

There are a few common choices for the mixing matrix $\tilde{\mathbf{W}}$, such as the Laplacian-based constant edge weight matrix [42] and the Metropolis constant edge weight matrix [54]. Let $\lambda_i(\tilde{\mathbf{W}})$ denote the $i$-th largest eigenvalue of $\tilde{\mathbf{W}}$ and $\sigma$ be the second largest magnitude eigenvalue of $\tilde{\mathbf{W}}$. Then, the following properties hold.

**Lemma 2.2.** *For $\tilde{\mathbf{W}}$ defined in Definition 2.1 and $\mathbf{W} := \tilde{\mathbf{W}} \otimes \mathbf{I}_p$, we have*

a. $1 = \lambda_1(\tilde{\mathbf{W}}) > \lambda_2(\tilde{\mathbf{W}}) \geq \ldots \geq \lambda_n(\tilde{\mathbf{W}}) > -1$;
b. $0 < \rho(\mathbf{W} - \mathbf{M}) = \sigma = \max\left\{|\lambda_2(\tilde{\mathbf{W}})|, |\lambda_n(\tilde{\mathbf{W}})|\right\} < 1$;
c. $\mathbf{M} = \mathbf{M}\mathbf{W} = \mathbf{W}\mathbf{M}$;
d. $\|\mathbf{W}^k \mathbf{x} - \mathbf{M}\mathbf{x}\| = \|(\mathbf{W}^k - \mathbf{M})(\mathbf{x} - \mathbf{M}\mathbf{x})\| \leq \sigma^k \|\mathbf{x} - \mathbf{M}\mathbf{x}\|$ *for any $\mathbf{x} \in \mathbb{R}^{np}$ and $k \geq 1$.*

*Proof.* The properties a-c holds, referred to [55]. Let us consider the property d. When $k = 1$, the inequality obviously holds. Then we show $(\mathbf{W} - \mathbf{M})^k = \mathbf{W}^k - \mathbf{M}$ for $k \geq 2$ by induction. For $k = 2$,

$$(\mathbf{W} - \mathbf{M})(\mathbf{W} - \mathbf{M})$$
$$= \mathbf{W}^2 - \mathbf{W}\mathbf{M} - \mathbf{M}\mathbf{W} + \mathbf{M} = \mathbf{W}^2 - \mathbf{M}.$$

For $k > 2$,

$$(\mathbf{W} - \mathbf{M})^k = (\mathbf{W} - \mathbf{M})^{k-1}(\mathbf{W} - \mathbf{M}) = (\mathbf{W}^{k-1} - \mathbf{M})(\mathbf{W} - \mathbf{M})$$
$$= \mathbf{W}^k - \mathbf{W}^{k-1}\mathbf{M} - \mathbf{M}\mathbf{W}^{k-1} + \mathbf{M} = \mathbf{W}^k - \mathbf{M}.$$

So the product of the matrix $\mathbf{W}$ converges to the average exponentially

$$\left\| \mathbf{W}^k - \mathbf{M} \right\| = \left\| (\mathbf{W} - \mathbf{M})^k \right\| \leq \left\| \mathbf{W} - \mathbf{M} \right\|^k \leq \sigma^k.$$

$\square$

## 2.1. Unified Nonconvex Algorithm Framework.

We propose a general framework of decentralized nonconvex methods, termed UDNA (Unified Decentralized Nonconvex Algorithm), as follows,

(2.2) $$\mathbf{x}^{t+1} = \mathbf{A}\mathbf{x}^t - \alpha\mathbf{B}\mathbf{H}^t\mathbf{v}^t,$$

(2.3) $$\mathbf{v}^{t+1} = \mathbf{C}\mathbf{v}^t + \mathbf{D}(\mathbf{g}^{t+1} - \mathbf{g}^t),$$

with initialization $\mathbf{v}^0 = \mathbf{g}^0$, where $\{\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \{\mathbf{H}^t\}\} \subset \mathbb{R}^{np \times np}$. $\mathbf{H}^t$ denotes a symmetric block diagonal matrix whose $i$-th block is $\mathbf{H}_i^t$, while $\mathbf{v}^t$ serves as a tracking variable that approximates to the global gradient $\nabla F(\bar{\mathbf{x}}^t)$. $\mathbf{A} = \tilde{\mathbf{A}} \otimes \mathbf{I}_p$, $\mathbf{B} = \tilde{\mathbf{B}} \otimes \mathbf{I}_p$, $\mathbf{C} = \tilde{\mathbf{C}} \otimes \mathbf{I}_p$, and $\mathbf{D} = \tilde{\mathbf{D}} \otimes \mathbf{I}_p$ are general symmetric matrices that satisfy the following structural assumption.

**Assumption 3.** *The matrices $\tilde{\mathbf{A}}$, $\tilde{\mathbf{B}}$, $\tilde{\mathbf{C}}$, and $\tilde{\mathbf{D}}$ are polynomial functions of the mixing matrix $\tilde{\mathbf{W}}$:*

(2.4) $$\tilde{\mathbf{A}} = \sum_{p=0}^{P_A} a_p \tilde{\mathbf{W}}^p, \quad \tilde{\mathbf{B}} = \sum_{p=0}^{P_B} b_p \tilde{\mathbf{W}}^p, \quad \tilde{\mathbf{C}} = \sum_{p=0}^{P_C} c_p \tilde{\mathbf{W}}^p, \text{ and } \tilde{\mathbf{D}} = \sum_{p=0}^{P_D} d_p \tilde{\mathbf{W}}^p,$$

*where $P_A, P_C \geq 1$ and $P_B, P_D \geq 0$. The coefficient sequences $\{a_p\}_{p=0}^{P_A}$, $\{b_p\}_{p=0}^{P_B}$, $\{c_p\}_{p=0}^{P_C}$, and $\{d_p\}_{p=0}^{P_D}$ are chosen such that $\tilde{\mathbf{A}}$ and $\tilde{\mathbf{C}}$ are doubly stochastic, while $\tilde{\mathbf{B}}$ and $\tilde{\mathbf{D}}$ are each either the identity matrix or a doubly stochastic matrix.*

By Assumption 3, the updating formula (2.3), and Lemma 2.2, an important property of the tracking variable $\mathbf{v}$ follows from induction that

(2.5) $$\mathbf{M}\mathbf{v}^t = \mathbf{M}\mathbf{g}^t \quad \Longleftrightarrow \quad \bar{\mathbf{v}}^t = \bar{\mathbf{g}}^t.$$

It is worth noting that the UDNA framework defined by (2.2) and (2.3) can subsume not only numerous state-of-the-art first-order methods but also various quasi-Newton methods, where $\mathbf{H}^t$ typically represents an approximation to the Hessian or its inverse. We now demonstrate the generality of the UDNA framework by showing how different choices of the matrices $\tilde{\mathbf{A}}, \tilde{\mathbf{B}}, \tilde{\mathbf{C}}, \tilde{\mathbf{D}}, \mathbf{H}^t$ recover various decentralized algorithmic schemes. Table 1 summarizes the specific choices of $\{\tilde{\mathbf{A}}, \tilde{\mathbf{B}}, \tilde{\mathbf{C}}, \tilde{\mathbf{D}}, \mathbf{H}^t\}$ in (2.2) and (2.3) yielding the desired equivalence.

At first, we consider the case $\mathbf{H}^t = \mathbf{I}$ for any $t$ to investigate the relationship with numerous gradient-based methods.

**ATC-GT method.** Taking $\tilde{\mathbf{A}} = \tilde{\mathbf{B}} = \tilde{\mathbf{C}} = \tilde{\mathbf{D}} = \tilde{\mathbf{W}}$ in (2.2) and (2.3), we obtain the recursion of the Adapt-Then-Combine Gradient-Tracking (ATC-GT) method as the following,

$$\mathbf{x}^{t+1} = \mathbf{W}(\mathbf{x}^t - \alpha\mathbf{v}^t),$$
$$\mathbf{v}^{t+1} = \mathbf{W}(\mathbf{v}^t + (\mathbf{g}^{t+1} - \mathbf{g}^t)),$$

which covers variants such as Aug-DGM [57] and ATC-DIGing [37].

**Non-ATC-GT method.** We notice the recursions of DIGing [36] and Harnessing [40] methods can be represented by (2.2) and (2.3) with $\tilde{\mathbf{A}} = \tilde{\mathbf{C}} = \tilde{\mathbf{W}}$ and $\tilde{\mathbf{B}} = \tilde{\mathbf{D}} = \mathbf{I}_p$, that is

$$\mathbf{x}^{t+1} = \mathbf{W}\mathbf{x}^t - \alpha\mathbf{v}^t,$$
$$\mathbf{v}^{t+1} = \mathbf{W}\mathbf{v}^t + (\mathbf{g}^{t+1} - \mathbf{g}^t).$$

These methods are of the non-ATC form and therefore classified as the non-ATC-GT method.

**Semi-ATC-GT method.** When substituting $\tilde{\mathbf{A}} = \tilde{\mathbf{B}} = \tilde{\mathbf{C}} = \tilde{\mathbf{W}}$ and $\tilde{\mathbf{D}} = \mathbf{I}_p$ into (2.2) and (2.3), we get the updates:

$$\mathbf{x}^{t+1} = \mathbf{W}(\mathbf{x}^t - \alpha\mathbf{v}^t),$$
$$\mathbf{v}^{t+1} = \mathbf{W}\mathbf{v}^t + (\mathbf{g}^{t+1} - \mathbf{g}^t),$$

which is equivalent to the recursion:

$$(2.6) \qquad \begin{aligned} \mathbf{x}^{t+2} =& \mathbf{W}\mathbf{x}^{t+1} - \alpha\mathbf{W}^2\mathbf{v}^t - \alpha\mathbf{W}(\mathbf{g}^{t+1} - \mathbf{g}^t) \\ =& 2\mathbf{W}\mathbf{x}^{t+1} - \mathbf{W}^2\mathbf{x}^t - \alpha\mathbf{W}(\mathbf{g}^{t+1} - \mathbf{g}^t). \end{aligned}$$

The equality (2.6) can be also obtained by eliminating the tracking variable $\mathbf{v}$ in (2.2) and (2.3) with $\tilde{\mathbf{A}} = \tilde{\mathbf{C}} = \tilde{\mathbf{D}} = \tilde{\mathbf{W}}$ and $\tilde{\mathbf{B}} = \mathbf{I}_p$. Related methods include SONATA [43] and NEXT [14], which can be classified as the semi-ATC-GT method.

Next, let us consider that the approximate matrix $\mathbf{H}^t$ is chosen as some special matrices. Several existing methods [46, 62, 18] leverage gradient-related variables to generate $\mathbf{H}^t$ as an approximation of the Hessian or its inverse matrix of the objective function. These methods indirectly exploit second-order information and are ususally classified as the quasi-Newton method, which are covered in our framework and listed below.

**DQN.** We can recover the DQN method [46] by setting $\tilde{\mathbf{A}} = \tilde{\mathbf{C}} = \tilde{\mathbf{D}} = \tilde{\mathbf{W}}$ and $\tilde{\mathbf{B}} = \tilde{\mathbf{W}}^2$ with $\mathbf{H}^t$ being the BFGS quasi-Newton matrix. It has the following recursion:

$$\mathbf{x}^{t+1} = \mathbf{W}(\mathbf{x}^t - \alpha\mathbf{d}^t),$$
$$\mathbf{v}^{t+1} = \mathbf{W}(\mathbf{v}^t + (\mathbf{g}^{t+1} - \mathbf{g}^t)),$$
$$\mathbf{d}^{t+1} = \mathbf{W}(\mathbf{H}^{t+1}\mathbf{v}^{t+1}),$$

where the $i$-th block of $\mathbf{H}^{t+1}$ is represented as

$$(2.7) \qquad \mathbf{H}_i^{t+1} = \left(\mathbf{I}_p - \frac{\mathbf{s}_i^t(\mathbf{y}_i^t)^\mathsf{T}}{(\mathbf{y}_i^t)^\mathsf{T}\mathbf{s}_i^t}\right)\mathbf{H}_i^t\left(\mathbf{I}_p - \frac{\mathbf{y}_i^t(\mathbf{s}_i^t)^\mathsf{T}}{(\mathbf{y}_i^t)^\mathsf{T}\mathbf{s}_i^t}\right) + \frac{\mathbf{s}_i^t(\mathbf{s}_i^t)^\mathsf{T}}{(\mathbf{y}_i^t)^\mathsf{T}\mathbf{s}_i^t}$$

with $\mathbf{s}_i^t = \mathbf{x}_i^{t+1} - \mathbf{x}_i^t$ and $\mathbf{y}_i^t = \mathbf{v}_i^{t+1} - \mathbf{v}_i^t$.

**DR-LM-DFP and D-LM-BFGS.** If we choose $\tilde{\mathbf{A}} = \tilde{\mathbf{C}} = \tilde{\mathbf{W}}$ and $\tilde{\mathbf{B}} = \tilde{\mathbf{D}} = \mathbf{I}_p$ in (2.2) and (2.3), and $\mathbf{H}^t$ is updated by damped regularized limited-memory DFP

or damped limited-memory BFGS techniques, we get the DR-LM-DFP or D-LM-BFGS methods [62].

**DGM-BB-C.** When setting $\tilde{\mathbf{A}} = \tilde{\mathbf{B}} = \tilde{\mathbf{C}} = \tilde{\mathbf{D}} = \tilde{\mathbf{W}}^K$ and constructing $\mathbf{H}^t$ as the BB block diagonal matrix, we recover the rescursion of the DGM-BB-C method [18] by substituting these into (2.2) and (2.3). Specifically, the updates become:

$$\mathbf{x}^{t+1} = \mathbf{W}^K(\mathbf{x}^t - \alpha\mathbf{H}^t\mathbf{v}^t),$$
$$\mathbf{v}^{t+1} = \mathbf{W}^K(\mathbf{v}^t + (\mathbf{g}^{t+1} - \mathbf{g}^t)).$$

The $i$-th block of $\mathbf{H}^{t+1}$ is updated by applying the BB method:

$$\mathbf{H}_i^{t+1} = \frac{\|\mathbf{s}_i^t\|^2}{(\mathbf{s}_i^t)^\mathsf{T}(\mathbf{g}_i^{t+1} - \mathbf{g}_i^t)}\mathbf{I}_p \text{ or } \frac{(\mathbf{s}_i^t)^\mathsf{T}(\mathbf{g}_i^{t+1} - \mathbf{g}_i^t)}{\|\mathbf{g}_i^{t+1} - \mathbf{g}_i^t\|^2}\mathbf{I}_p.$$

**DSG.** Suppose $\tilde{\mathbf{A}} = \tilde{\mathbf{C}} = \tilde{\mathbf{W}}$ and $\tilde{\mathbf{B}} = \tilde{\mathbf{D}} = \mathbf{I}_p$. We consider $\mathbf{H}^{t+1}$ behaves as a spectral-like diagonal matrix with the $i$-th block $\mathbf{H}_i^{t+1}=(\delta_i^{t+1})^{-1}\mathbf{I}_p$, where

$$(2.8) \quad \delta_i^{t+1} = \text{Proj}_{[\delta_{\min}, \delta_{\max}]}\left\{\frac{(\mathbf{s}_i^t)^\mathsf{T}(\mathbf{g}_i^{t+1} - \mathbf{g}_i^t)}{\|\mathbf{s}_i^t\|^2} + \delta_i^t\sum_{j\in\mathcal{N}_i}\tilde{W}_{ij}\left(1 - \frac{(\mathbf{s}_j^t)^\mathsf{T}\mathbf{s}_i^t}{\|\mathbf{s}_i^t\|^2}\right)\right\},$$

and $0 < \delta_{\min} < \delta_{\max} < \infty$ are parameters. Then we obtain the DSG method [25].

We formally present the detailed UDNA algorithmic procedure in Algorithm 1, and all the aforementioned methods can follow this procedure.

---

**Algorithm 1** Unified Decentralized Nonconvex Algorithm (UDNA)

---

**Input:** Initial point $\mathbf{x}^0$, Maximum iteration T, Stepsize $\alpha > 0$, Parameters $\tilde{\mathbf{A}}$, $\tilde{\mathbf{B}}$, $\tilde{\mathbf{C}}$, $\tilde{\mathbf{D}}$.
 1: Set $t = 0$ and $\mathbf{v}^0 = \mathbf{g}^0$. Initialize $\mathbf{H}^0$.
 2: If $t \geq T$, stop.
 3: $\mathbf{x}^{t+1} = \mathbf{A}\mathbf{x}^t - \alpha\mathbf{B}\mathbf{H}^t\mathbf{v}^t$.
 4: $\mathbf{v}^{t+1} = \mathbf{C}\mathbf{v}^t + \mathbf{D}(\mathbf{g}^{t+1} - \mathbf{g}^t)$.
 5: Update $\mathbf{H}^{t+1}$.
 6: Set $t = t + 1$ and go to Step 2.
**Output:** $\mathbf{x}^T$.

---

*Remark* 2.3. (Relationship with other frameworks). It should be noticed that Berahas et al. [8] proposed a Gradient Tracking Algorithmic (GTA) framework which unifies various gradient tracking methods. All the methods covered in the GTA framework are also included in our framework. Moreover, the convergence of GTA is only studied for the strongly convex objectives. In contrast, we analyse the convergence under the general nonconvex settings.

A stochastic nonconvex unified framework, termed SUDA [4], was previously proposed, integrating several well-known decentralized methods within a primal-dual formulation. While both SUDA and our UDNA encompass gradient tracking methods, a key distinction of our work lies in the integration of quasi-Newton methods. Our framework successfully addresses the convergence guarantees for these quasi-Newton variants applied in nonconvex settings, which remains a challenging and less explored area. Moreover, we employ a different analytical framework of subsequence convergence compared to SUDA and, more importantly, establish the

TABLE 1. Special cases in our framework

| Methods | $\check{\mathbf{A}}$ | $\check{\mathbf{B}}$ | $\check{\mathbf{C}}$ | $\check{\mathbf{D}}$ | $\mathbf{H}^t$ |
|---|---|---|---|---|---|
| DIGing [36], Harnessing [40] | $\tilde{\mathbf{W}}$ | $\mathbf{I}_p$ | $\tilde{\mathbf{W}}$ | $\mathbf{I}_p$ | $\mathbf{I}$ |
| SONATA [43], NEXT [14] | $\tilde{\mathbf{W}}$ | $\tilde{\mathbf{W}}(\mathbf{I}_p)$ | $\tilde{\mathbf{W}}$ | $\mathbf{I}_p(\tilde{\mathbf{W}})$ | $\mathbf{I}$ |
| Aug-DGM [57], ATC-DIGing [37] | $\tilde{\mathbf{W}}$ | $\tilde{\mathbf{W}}$ | $\tilde{\mathbf{W}}$ | $\tilde{\mathbf{W}}$ | $\mathbf{I}$ |
| DR-LM-DFP [62], D-LM-BFGS [62] | $\tilde{\mathbf{W}}$ | $\mathbf{I}_p$ | $\tilde{\mathbf{W}}$ | $\mathbf{I}_p$ | BFGS/DFP ‡ |
| DQN [47] | $\tilde{\mathbf{W}}$ | $\tilde{\mathbf{W}}^2$ | $\tilde{\mathbf{W}}$ | $\tilde{\mathbf{W}}$ | BFGS |
| DGM-BB-C [18] | $\tilde{\mathbf{W}}^{K†}$ | $\tilde{\mathbf{W}}^K$ | $\tilde{\mathbf{W}}^K$ | $\tilde{\mathbf{W}}^K$ | BB |
| DSG[25] | $\tilde{\mathbf{W}}$ | $\mathbf{I}_p$ | $\tilde{\mathbf{W}}$ | $\mathbf{I}_p$ | (2.8) |

† Taking the mixing matrix $\tilde{\mathbf{W}}$ to the $K$ power represents performing $K$ communications at one iteration.

‡ "BFGS", "DFP", and "BB" means the approximate matrix $\mathbf{H}^t$ generated by the BFGS, DFP, or BB quasi-Newton formular.

TABLE 2. Comparisons with existing frameworks

| Frameworks | GTA[8] | ABC[56] | SUDA [4] | UDNA(Ours) |
|---|---|---|---|---|
| Nonconvex or not | No | No | Yes | Yes |
| Cover GT? | Yes | Yes | Yes | Yes |
| Cover ED? | No | Yes | Yes | No |
| Cover QN? | No | No | No | Yes |
| Use KL? | No | No | No | Yes |

stronger result of whole sequence convergence by leveraging the KL property of the objective function.

Table 2 summarizes the comparsions with several existing frameworks, where GT, ED, and QN respectively represent gradient tracking, exact diffusion, and quasi-Newton methods. The ABC framework [56], while built on the similar primal-dual idea to SUDA[4], is restricted to strongly convex objectives.

2.2. **Global convergence.** We now analyze the global convergence of the proposed unified algorithm UDNA for minimizing (1.1) in which the local objective function $f_i$, $i = 1, \ldots, n$, is Lipschitz continuously differentiable, but possibly nonconvex. For notational convenience, we define

$$\sigma_A = \rho(\mathbf{A} - \mathbf{M}), \quad \sigma_B = \rho(\mathbf{B} - \mathbf{M}),$$
$$\sigma_C = \rho(\mathbf{C} - \mathbf{M}), \quad \sigma_D = \rho(\mathbf{D} - \mathbf{M}).$$

The convergence analysis relies on the following fundamental assumption, which require the sequence of approximate matrices $\{\mathbf{H}_i^t\}$ to have uniformly bounded eigenvalues in the subspace spanned by $\{\mathbf{v}_i^t\}$. This provides the necessary control over the approximation quality of the Hessian or its inverse, allowing us to establish a unified convergence theory for all algorithms subsumed by the UDNA framework, regardless of the specific choice of $\mathbf{H}_i^t$.

**Assumption 4.** *The approximate matrices $\{\mathbf{H}_i^t\}$ satisfy*

$$\psi\|\mathbf{v}_i^t\|^2 \leq (\mathbf{v}_i^t)^\mathsf{T}\mathbf{H}_i^t\mathbf{v}_i^t, \text{ and } \|\mathbf{H}_i^t\mathbf{v}_i^t\|^2 \leq \Psi^2\|\mathbf{v}_i^t\|^2,$$

*for any $t$ and $i$, where $\Psi \geq \psi > 0$ and $\{\mathbf{v}_i^t\}$ are from (2.3).*

This assumption implies that the eigenvalues of $\mathbf{H}_i^t$ in the direction of $\mathbf{v}_i^t$ lie within the interval $[\psi, \Psi]$. However, ensuring eigenvalue boundedness of $\mathbf{H}^t$ is a major challenge in decentralized nonconvex settings. To address this critical issue, we will provide several schemes different from the previously proposed ones in [62]. We will not describe how to construct $\mathbf{H}^t$ until in Subsection 2.4, where specific updating schemes for $\mathbf{H}^t$ satisfying Assumption 4 will be proposed. Recall the notation $\overline{\nabla} f(\mathbf{x}^t) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_i^t)$. Then, we obtain $\mathbf{1} \otimes \overline{\nabla} f(\mathbf{x}^t) = \mathbf{M} \mathbf{g}^t$. Note that if $\mathbf{x}^* \in \mathbb{R}^{np}$ satisfies

$$\|\overline{\nabla} f(\mathbf{x}^*)\|^2 + \|\mathbf{x}^* - \mathbf{M}\mathbf{x}^*\|^2 = 0, \tag{2.9}$$

then we have $\mathbf{x}_1^* = \ldots = \mathbf{x}_i^* =: \mathbf{z}^*$ and $\mathbf{z}^*$ would be a first-order stationary point of (1.1). To obtain (2.9), it suffices to show $\|\mathbf{v}^t\| \to 0$, $\|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\| \to 0$, and $\|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\| \to 0$ as $t \to \infty$ since $\mathbf{M}\mathbf{v}^t = \mathbf{M}\mathbf{g}^t$ for any $t$ by (2.5). Hence, by (2.9), we say $\mathbf{x}^t$ is an $\epsilon$-stationary solution for some $\epsilon > 0$ if

$$\|\mathbf{v}^t\|^2 + \|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2 + \|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2 \leq \epsilon. \tag{2.10}$$

Given any $\epsilon > 0$, to show UDNA will generate a $(\mathbf{x}^t, \mathbf{v}^t)$ satisfying (2.10), we define the following potential function

$$P(\mathbf{x}^t, \mathbf{v}^t) = F(\bar{\mathbf{x}}^t) + \|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2 + \frac{1 - (1+\tau)\sigma_A^2}{4(1 + 1/\eta)L^2\sigma_D^2}\|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2, \tag{2.11}$$

where $\eta$ and $\tau$ are some positive constants. The term $F(\bar{\mathbf{x}}^t)$ corresponds to the objective value at the average of local variables, while $\|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2$ measures the consensus error, i.e., the deviation of the local variables from their global average. The term $\|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2$ captures the gradient tracking error, representing the difference between the average-gradient approximations and the true average gradient. Since each $f_i$ is lower bounded, $P$ is also bounded below. To establish the sufficient descent property of $P(\mathbf{x}^t, \mathbf{v}^t)$, we begin by deriving a recursion relationship of $F(\bar{\mathbf{x}}^t)$.

**Lemma 2.4.** *Suppose that Assumptions 1, 2, 3, and 4 hold. Let $\{\mathbf{x}^t\}$ be the sequence generated by Algorithm 1 (UDNA). We have for all $t \geq 0$,*

$$F(\bar{\mathbf{x}}^{t+1}) \leq F(\bar{\mathbf{x}}^t) - \left(\frac{\alpha\psi}{2n} - \frac{L\alpha^2\Psi^2}{2n}\right)\|\mathbf{v}^t\|^2 \tag{2.12}$$

$$+ \frac{L^2\alpha\Psi^2}{n}\|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2 + \frac{\alpha\Psi^2}{n}\|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2.$$

*Proof.* By the $L$-Lipschitz continuity of $\nabla F$, we have

$$F(\bar{\mathbf{x}}^{t+1}) \leq F(\bar{\mathbf{x}}^t) + \left\langle \nabla F(\bar{\mathbf{x}}^t), \bar{\mathbf{x}}^{t+1} - \bar{\mathbf{x}}^t \right\rangle + \frac{L}{2}\|\bar{\mathbf{x}}^{t+1} - \bar{\mathbf{x}}^t\|^2, \tag{2.13}$$

Define an auxiliary sequence $\{\tilde{\mathbf{x}}^t\}_{t=0}^T$ such that for each node $i$, $\tilde{\mathbf{x}}_i^{t+1} = \mathbf{x}_i^t - \alpha\mathbf{H}_i^t\mathbf{v}_i^t$ and $\tilde{\mathbf{x}}_i^0 = \mathbf{x}_i^0$, and observe that $\frac{1}{n}\sum_{i=1}^n \tilde{\mathbf{x}}_i^t = \frac{1}{n}\sum_{i=1}^n \mathbf{x}_i^t = \bar{\mathbf{x}}^t$. We now decompose

the inner product term in (2.13) as follows:

$$
(2.14) \qquad \left\langle \nabla F(\bar{\mathbf{x}}^t), \bar{\mathbf{x}}^{t+1} - \bar{\mathbf{x}}^t \right\rangle = \frac{1}{n} \sum_{i=1}^n \left\langle \nabla F(\bar{\mathbf{x}}^t), \tilde{\mathbf{x}}_i^{t+1} - \mathbf{x}_i^t \right\rangle
$$

$$
= \underbrace{\frac{1}{n} \sum_{i=1}^n \left\langle \mathbf{v}_i^t, \tilde{\mathbf{x}}_i^{t+1} - \mathbf{x}_i^t \right\rangle}_{\text{term (I)}} + \underbrace{\frac{1}{n} \sum_{i=1}^n \left\langle \nabla F(\bar{\mathbf{x}}^t) - \overline{\nabla} f(\mathbf{x}^t), \tilde{\mathbf{x}}_i^{t+1} - \mathbf{x}_i^t \right\rangle}_{\text{term (II)}}
$$

$$
+ \underbrace{\frac{1}{n} \sum_{i=1}^n \left\langle \overline{\nabla} f(\mathbf{x}^t) - \mathbf{v}_i^t, \tilde{\mathbf{x}}_i^{t+1} - \mathbf{x}_i^t \right\rangle}_{\text{term (III)}}.
$$

We now bound each term separately.

**Term (I):** By Assumption 4, we have

$$
(2.15) \qquad\qquad \frac{1}{n} \sum_{i=1}^n \left\langle \mathbf{v}_i^t, \tilde{\mathbf{x}}_i^{t+1} - \mathbf{x}_i^t \right\rangle \leq -\frac{\alpha\psi}{n} \|\mathbf{v}^t\|^2.
$$

**Term (II):** Using Young's inequality with some $c > 0$ and Assumption 4, we deduce

$$
(2.16) \quad \frac{1}{n} \sum_{i=1}^n \left\langle \nabla F(\bar{\mathbf{x}}^t) - \overline{\nabla} f(\mathbf{x}^t), \tilde{\mathbf{x}}_i^{t+1} - \mathbf{x}_i^t \right\rangle
$$

$$
\leq \frac{c}{2} \|\nabla F(\bar{\mathbf{x}}^t) - \overline{\nabla} f(\mathbf{x}^t)\|^2 + \frac{\alpha^2 \Psi^2}{2nc} \|\mathbf{v}^t\|^2 \leq \frac{L^2 c}{2n} \|\mathbf{M}\mathbf{x}^t - \mathbf{x}^t\|^2 + \frac{\alpha^2 \Psi^2}{2nc} \|\mathbf{v}^t\|^2,
$$

where the last inequality follows from the $L$-Lipschitz continuity of $\nabla f_i$ and the relation $\mathbf{1} \otimes \bar{\mathbf{x}}^t = \mathbf{M}\mathbf{x}^t$.

**Term (III):** Applying Young's inequality with some $d > 0$, Assumption 4, and the relation $\mathbf{1} \otimes \overline{\nabla} f(\mathbf{x}^t) = \mathbf{M}\mathbf{v}^t$ from (2.5), we derive

$$
(2.17) \quad \frac{1}{n} \sum_{i=1}^n \left\langle \overline{\nabla} f(\mathbf{x}^t) - \mathbf{v}_i^t, \tilde{\mathbf{x}}_i^{t+1} - \mathbf{x}_i^t \right\rangle
$$

$$
\leq \frac{d}{2n} \|\mathbf{1} \otimes \overline{\nabla} f(\mathbf{x}^t) - \mathbf{v}^t\|^2 + \frac{\alpha^2 \Psi^2}{2nd} \|\mathbf{v}^t\|^2 \leq \frac{d}{2n} \|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2 + \frac{\alpha^2 \Psi^2}{2nd} \|\mathbf{v}^t\|^2.
$$

Plugging (2.15), (2.16), and (2.17) back in (2.14) yields

$$
(2.18) \quad \left\langle \nabla F(\bar{\mathbf{x}}^t), \bar{\mathbf{x}}^{t+1} - \bar{\mathbf{x}}^t \right\rangle
$$

$$
\leq -\left( \frac{\alpha\psi}{n} - \frac{\alpha^2 \Psi^2}{2n} \left( \frac{1}{c} + \frac{1}{d} \right) \right) \|\mathbf{v}^t\|^2 + \frac{L^2 c}{2n} \|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2 + \frac{d}{2n} \|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2.
$$

For term $\|\bar{\mathbf{x}}^{t+1} - \bar{\mathbf{x}}^t\|^2$ in (2.13), it follows from Assumption 4 that

$$
(2.19) \qquad\qquad \|\bar{\mathbf{x}}^{t+1} - \bar{\mathbf{x}}^t\|^2 = \frac{1}{n} \|\mathbf{M}(\tilde{\mathbf{x}}^{t+1} - \mathbf{x}^t)\|^2 \leq \frac{\alpha^2 \Psi^2}{n} \|\mathbf{v}^t\|^2
$$

Finally, substituting (2.18) and (2.19) into (2.13), and choosing $c = 2\alpha\Psi^2/\psi$, $d = 2\alpha\Psi^2/\psi$, we get obtain the desired inequality (2.12). $\qquad\square$

**Theorem 2.5.** *Suppose that Assumptions 1, 2, 3, and 4 hold. Let $\{\mathbf{x}^t\}$ be the sequence generated by Algorithm 1 (UDNA). If*

$$(2.20) \qquad \alpha \leq \min\left\{ \frac{(1-\sigma_A^2)(1-\sigma_C^2)^2 n}{64L^2\sigma_D^2\Psi^2}, \frac{(1-\sigma_A^2)\psi}{(2L+n+8\sigma_B^2 n)\Psi^2}, \frac{(1-\sigma_A^2)n}{2L^2\Psi^2} \right\},$$

*then we have the following convergence rate for UDNA*

$$(2.21) \quad \min_{0 \leq t \leq T}\left\{ \|\mathbf{v}^t\|^2 + \|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2 + \|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2 \right\} \leq \frac{P(\mathbf{x}^0, \mathbf{v}^0) - P(\mathbf{x}^T, \mathbf{v}^T)}{\gamma \min\{\alpha, 1\}T},$$

*where $P$ is the potential function defined in (2.11) with $\tau = \frac{1-\sigma_A^2}{2\sigma_A^2}$ and $\eta = \frac{1-\sigma_C^2}{2\sigma_C^2}$, and $\gamma = \min\{a_1, a_2, a_3\}$ with*

$$(2.22) \qquad \begin{cases} a_1 = \frac{\psi}{2n} - \frac{(2L+n+8\sigma_B^2 n)\Psi^2}{4n(1-\sigma_A^2)}\alpha \geq \frac{\psi}{4n} > 0, \\ a_2 = \frac{(1-\sigma_A^2)(1-\sigma_C^2)^2}{32L^2\sigma_D^2} - \frac{\Psi^2\alpha}{n} \geq \frac{(1-\sigma_A^2)(1-\sigma_C^2)^2}{64L^2\sigma_D^2}, \quad and \\ a_3 = 1 - \sigma_A^2 - \frac{L^2\Psi^2\alpha}{n} \geq \frac{1-\sigma_A^2}{2} > 0. \end{cases}$$

*Proof.* At first, we establish a recursive upper bound for $\|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2$. From the iteration (2.2), we have

$$(2.23) \qquad \|\mathbf{x}^{t+1} - \mathbf{M}\mathbf{x}^{t+1}\|^2 = \|\mathbf{A}\mathbf{x}^t - \mathbf{M}\mathbf{x}^t - \alpha\mathbf{B}\mathbf{H}^t\mathbf{v}^t + \alpha\mathbf{M}\mathbf{H}^t\mathbf{v}^t\|^2$$

$$\leq (1+\tau)\|\mathbf{A}\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2 + (1+1/\tau)\alpha^2\|(\mathbf{B}-\mathbf{M})\mathbf{H}^t\mathbf{v}^t\|^2$$

$$\leq (1+\tau)\sigma_A^2\|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2 + (1+1/\tau)\alpha^2\sigma_B^2\Psi^2\|\mathbf{v}^t\|^2$$

where the first inequality applies Young's inequality with some $\tau > 0$ and the second inequality uses Lemma 2.2 and Assumption 4.

Next, we establish a recursive upper bound for $\|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2$. From the updating formluar (2.3), we have

$$\|\mathbf{v}^{t+1} - \mathbf{M}\mathbf{v}^{t+1}\|^2 = \|\mathbf{C}\mathbf{v}^t + \mathbf{D}\mathbf{g}^{t+1} - \mathbf{D}\mathbf{g}^t - \mathbf{M}\mathbf{v}^t - \mathbf{M}\mathbf{g}^{t+1} + \mathbf{M}\mathbf{g}^t\|^2$$

$$\leq (1+\eta)\|\mathbf{C}\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2 + (1+1/\eta)\|(\mathbf{D}-\mathbf{M})(\mathbf{g}^{t+1}-\mathbf{g}^t)\|^2$$

$$\leq (1+\eta)\sigma_C^2\|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2 + (1+1/\eta)\sigma_D^2 L^2\|\mathbf{x}^{t+1} - \mathbf{x}^t\|^2,$$

where the first inequality follows from Young's inequality with some $\eta > 0$ and the second inequality uses Lemma 2.2 and the $L$-Lipschitz continuity of $\nabla f_i$. Note that $\mathbf{x}^{t+1} - \mathbf{x}^t = (\mathbf{A} - \mathbf{I})(\mathbf{x}^t - \mathbf{M}\mathbf{x}^t) - \alpha\mathbf{B}\mathbf{H}^t\mathbf{v}^t$, and $\rho(\mathbf{A} - \mathbf{I}) < 2$. Then, invoking Lemma A.2 and Assumption 4 yields

$$(2.24) \qquad \|\mathbf{v}^{t+1} - \mathbf{M}\mathbf{v}^{t+1}\|^2 \leq \|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2 - (1 - (1+\eta)\sigma_C^2)\|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2$$

$$+ 8(1+1/\eta)\sigma_D^2 L^2\|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2 + 2(1+1/\eta)\sigma_D^2 L^2\alpha^2\Psi^2\|\mathbf{v}^t\|^2.$$

Multiplying both sides of the above inequality by $\frac{1-(1+\tau)\sigma_A^2}{4(1+1/\eta)L^2\sigma_D^2}$, we have

$$(2.25)$$

$$\frac{1-(1+\tau)\sigma_A^2}{4(1+1/\eta)L^2\sigma_D^2}\|\mathbf{v}^{t+1} - \mathbf{M}\mathbf{v}^{t+1}\|^2$$

$$\leq \frac{1-(1+\tau)\sigma_A^2}{4(1+1/\eta)L^2\sigma_D^2}\|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2 - \frac{(1-(1+\tau)\sigma_A^2)(1-(1+\eta)\sigma_C^2)}{4(1+1/\eta)L^2\sigma_D^2}\|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2$$

$$+ 2(1-(1+\tau)\sigma_A^2)\|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2 + \frac{(1-(1+\tau)\sigma_A^2)\alpha^2\Psi^2}{2}\|\mathbf{v}^t\|^2.$$

Adding up (2.12), (2.25), and (2.23) yields

$$P(\mathbf{x}^{t+1}, \mathbf{v}^{t+1}) \leq P(\mathbf{x}^t, \mathbf{v}^t) - \left(2(1 - (1+\tau)\sigma_A^2) - \frac{L^2\alpha\Psi^2}{n}\right)\|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2$$

$$- \left(\frac{\alpha\psi}{2n} - \frac{L\alpha^2\Psi^2}{2n} - \frac{(1-(1+\tau)\sigma_A^2)\alpha^2\Psi^2}{2} - \left(1 + \frac{1}{\tau}\right)\sigma_B^2\alpha^2\Psi^2\right)\|\mathbf{v}^t\|^2$$

$$- \left(\frac{(1-(1+\tau)\sigma_A^2)(1-(1+\eta)\sigma_C^2)}{4(1+1/\eta)L^2\sigma_D^2} - \frac{\alpha\Psi^2}{n}\right)\|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2.$$

Setting $\tau = \frac{1-\sigma_A^2}{2\sigma_A^2}$ and $\eta = \frac{1-\sigma_C^2}{2\sigma_C^2}$, and and using the fact that $0 < \sigma < 1$, we derive

$$P(\mathbf{x}^{t+1}, \mathbf{v}^{t+1}) \leq P(\mathbf{x}^t, \mathbf{v}^t) - \left(1 - \sigma_A^2 - \frac{L^2\Psi^2\alpha}{n}\right)\|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2$$

$$- \left(\frac{\psi}{2n} - \frac{L\Psi^2\alpha}{2n} - \frac{(1-\sigma_A^2)\Psi^2\alpha}{4} - \frac{(1+\sigma_A^2)\sigma_B^2\Psi^2\alpha}{1-\sigma_A^2}\right)\alpha\|\mathbf{v}^t\|^2$$

$$- \left(\frac{(1-\sigma_A^2)(1-\sigma_C^2)^2}{16(1+\sigma_C^2)L^2\sigma_D^2} - \frac{\Psi^2\alpha}{n}\right)\|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2$$

$$\leq P(\mathbf{x}^t, \mathbf{v}^t) - \left(1 - \sigma_A^2 - \frac{L^2\Psi^2\alpha}{n}\right)\|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2$$

$$- \left(\frac{\psi}{2n} - \frac{(2L + n + 8\sigma_B^2 n)\Psi^2}{4n(1-\sigma_A^2)}\alpha\right)\alpha\|\mathbf{v}^t\|^2$$

$$- \left(\frac{(1-\sigma_A^2)(1-\sigma_C^2)^2}{32L^2\sigma_D^2} - \frac{\Psi^2\alpha}{n}\right)\|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2.$$

Hence, by the definitions of $a_1$, $a_2$, and $a_3$ in (2.22), we obtain

(2.26) $P(\mathbf{x}^{t+1}, \mathbf{v}^{t+1}) - P(\mathbf{x}^t, \mathbf{v}^t) \leq -a_1\alpha\|\mathbf{v}^t\|^2 - a_2\|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2 - a_3\|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2,$

where by direct calculation we have

$$a_1 \geq \frac{\psi}{4n}, \quad a_2 \geq \frac{(1-\sigma_A^2)(1-\sigma_C^2)^2}{64L^2\sigma_D^2}, \quad \text{and} \quad a_3 \geq \frac{1-\sigma^2}{2}.$$

Summing (2.26) over $t = 0, \ldots, T$ and dividing both sides by $\gamma \min\{\alpha, 1\}$, we get the desired convergence rate (2.21), where $\gamma = \min\{a_1, a_2, a_3\}$. $\square$

Given any $\epsilon > 0$, by Theorem 2.5, UDNA will take at most $P(\mathbf{x}^0, \mathbf{v}^0)(\epsilon\gamma \min\{\alpha, 1\})^{-1}$ iterations to generate an $\epsilon$-stationary point $\mathbf{x}^t$ satisfying (2.10), where $\gamma$ is the constant given in (2.21).

**Corollary 2.6.** *Under the assumptions of Theorem 2.5 and the choice of stepsize $\alpha$ given by (2.20), it holds for the sequence $\{\mathbf{x}^t\}$ generated by Algorithm 1 (UDNA) that $\{\mathbf{x}^t\}$ has at least one accumulation point and any such point is a stationary point of the problem (1.1).*

*Proof.* Since the sequences $\{\|\mathbf{v}^t\|^2\}$, $\{\|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|^2\}$, and $\{\|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\|^2\}$ are summable, we obtain

(2.27) $\|\mathbf{v}^t\|, \|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|, \|\mathbf{v}^t - \mathbf{M}\mathbf{v}^t\| \to 0 \Rightarrow \|\nabla F(\bar{\mathbf{x}}^t)\| \to 0, F(\bar{\mathbf{x}}^t) \to P^\infty,$

for some $P^\infty \in \mathbb{R}$. We now show the sequence $\mathbf{x}^t$ is bounded. By Theorem 2.5 and the stepsize choice (2.20), the potential function $P$ is nonincreasing and bounded above by $P(\mathbf{x}^0, \mathbf{v}^0) < \infty$. Hence, $F(\bar{\mathbf{x}}^t) \leq P(\mathbf{x}^0, \mathbf{v}^0)$ indicates that $\bar{\mathbf{x}}^t$ is bounded

due to the coercivity of $F$; that is, there exists some $B_1 > 0$ such that $\|\bar{\mathbf{x}}^t\| \leq B_1$ for any $t$. Since $\|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\| \to 0$, $\mathbf{x}^t - \mathbf{M}\mathbf{x}^t$ is also bounded: $\|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\| \leq B_2$ for some $B_2 > 0$ and all $t$. Consequently,

$$\|\mathbf{x}^t\| \leq \|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\| + \sqrt{n}\|\bar{\mathbf{x}}^t\| \leq \sqrt{n}B_1 + B_2, \ \forall \ t \geq 0.$$

The boundedness of $\mathbf{x}^t$ guarantees that there exist a convergent subsequence and a limit point, denoted by $\{\mathbf{x}^{t_k}\}_{k \in \mathbb{N}} \to \mathbf{x}^\infty$ as $k \to \infty$. From (2.27), we have $\|\nabla F(\bar{\mathbf{x}}^{t_k})\| + \|\mathbf{x}^{t_k} - \mathbf{M}\mathbf{x}^{t_k}\| \to 0$ as $k \to \infty$, which implies $\|\nabla F(\bar{\mathbf{x}}^\infty)\| + \|\mathbf{x}^\infty - \mathbf{M}\mathbf{x}^\infty\| = 0$. $\qquad\square$

While **Theorem 2.5** certifies stationarity of every limit point of $\{\mathbf{x}^t\}$ generated by UDNA, the convergence of the whole sequences is not guaranteed. The next subsection addresses this issue given that the aggregated objective function $F$ has the KŁ property.

### 2.3. Convergence under KŁ property.
We now strengthen the subsequence convergence result in Theorem 2.5, proving the global convergence of the sequence $\{\mathbf{x}^t\}$ by postulating that $F$ has the KŁ property. Notably, the convergence rate naturally emerges as a direct consequence of this analysis.

**Definition 2.7.** (KŁ property & KŁ function) A proper closed function $h$ is said to have the Kurdyka-Łojasiewicz (KŁ) property at $\mathbf{z} \in \operatorname{dom} h$ if there exist a neighborhood $\mathcal{N}$ of $\mathbf{z}$, $v \in (0, \infty]$ and a continuous concave function $\phi : [0, v) \to \mathbb{R}_+$ with $\phi(0) = 0$ such that

1. $\phi$ is continuously differentiable on $(0, v)$ with $\phi' > 0$.
2. For all $\tilde{\mathbf{z}} \in \mathcal{N}$ with $h(\mathbf{z}) < h(\tilde{\mathbf{z}}) < h(\mathbf{z}) + v$, it holds that

$$\phi'(h(\tilde{\mathbf{z}}) - h(\mathbf{z})) \operatorname{dist}(0, \partial h(\tilde{\mathbf{z}})) \geq 1.$$

A proper closed function $h$ satisfying the KŁ property at all points in $\operatorname{dom} h$ is called a KŁ function.

Of particular interest in our analysis is the notion of the KŁ exponent, defined as follows.

**Definition 2.8.** (KŁ exponent) For a proper closed function $h$ satisfying the KŁ property at $\mathbf{z} \in \operatorname{dom} h$, if the desingularizing function $\phi$ can be chosen as $\phi(s) = \kappa s^{1-\theta}$ for some $\kappa > 0$ and $\theta \in [0, 1)$, i.e., there exists $\epsilon > 0$ and $v \in (0, \infty]$ such that

$$\operatorname{dist}(0, \partial h(\tilde{\mathbf{z}})) \geq \kappa |h(\tilde{\mathbf{z}}) - h(\mathbf{z})|^\theta$$

whenever $\|\tilde{\mathbf{z}} - \mathbf{z}\| \leq \epsilon$ and $0 < |h(\tilde{\mathbf{z}}) - h(\mathbf{z})| < v$, then $h$ is said to have the KŁ property at $\mathbf{z}$ with exponent $\theta$. If $h$ is a KŁ function and has the same exponent $\theta$ in $\operatorname{dom} h$, then we say that $h$ is a KŁ function with exponent $\theta$.

Throughout this subsection, we introduce the following notations to facilitate our theoretical analysis.

$$(2.28) \quad \mathbf{X} = [\mathbf{x}_1^\mathsf{T}; \mathbf{x}_2^\mathsf{T}; \ldots; \mathbf{x}_n^\mathsf{T}] \in \mathbb{R}^{n \times p}, \ \mathbf{V} = [\mathbf{v}_1^\mathsf{T}; \mathbf{v}_2^\mathsf{T}; \ldots; \mathbf{v}_n^\mathsf{T}] \in \mathbb{R}^{n \times p}, \ \tilde{\mathbf{M}} = \frac{1}{n}\mathbf{1}_n^\mathsf{T}\mathbf{1}_n.$$

Then, $F(\bar{\mathbf{x}}^t) = F((1/n)\mathbf{1}_n^\mathsf{T}\mathbf{X}^t) := G(\mathbf{X}^t)$. The following lemma shows that the KŁ property of $F$ naturally transfers to $G$.

**Lemma 2.9.** *Under the same assumptions and parameter choices of Theorem 2.5, let $\mathbf{1}_n \otimes (\bar{\mathbf{x}}^\infty)^\mathsf{T} = \mathbf{X}^\infty$ be an accumulation point of $\mathbf{X}^t$ where $\bar{\mathbf{x}}^\infty$ is some critical point of $F$. If $F$ has the KL property at $\bar{\mathbf{x}}^\infty$ with exponent $\theta_F \in [0, 1)$ and parameters $\{\kappa_F > 0, v_F > 0\}$, then $G$ inherits the KL property at $\mathbf{X}^\infty$ with exponent $\theta_G = \theta_F$ and parameters $\{\kappa_G = \frac{\kappa_F}{\sqrt{n}}, v_G = v_F\}$.*

*Proof.* By the KL property of $F$, there exists a neighborhood $V = \{\mathbf{z} \in \mathbb{R}^{n \times p} : \|\mathbf{z} - \bar{\mathbf{x}}^\infty\| \leq \epsilon_F\}$ of $\bar{\mathbf{x}}^\infty$, a constant $\kappa_F > 0$, and an exponent $\theta_F \in [0, 1)$ such that for all $\mathbf{z} \in V \cap \{\mathbf{z} : 0 < |F(\mathbf{z}) - F(\bar{\mathbf{x}}^\infty)| < v\}$, the following inequality holds,

$$\|\nabla F(\mathbf{z})\| \geq \kappa_F |F(\mathbf{z}) - F(\bar{\mathbf{x}}^\infty)|^{\theta_F}.$$

We define a neighborhood $U$ of $\mathbf{X}^\infty$:

$$U = \{\mathbf{X} \in \mathbb{R}^{n \times p} : \|\mathbf{X} - \mathbf{X}^\infty\| \leq \epsilon_G\}.$$

We can choose $\epsilon_G$ to be sufficiently small (specifically, $\epsilon_G = \sqrt{n}\epsilon_F$) such that for any $\mathbf{X} \in U$, the corresponding average $\bar{\mathbf{x}} = \frac{1}{n}\mathbf{1}^\mathsf{T}\mathbf{X}$ lies within the neighborhood $V$ of $\bar{\mathbf{x}}^\infty$. This is guaranteed because $\|\bar{\mathbf{x}} - \bar{\mathbf{x}}^\infty\| \leq \frac{1}{\sqrt{n}}\|\mathbf{X} - \mathbf{X}^\infty\|$. Note $\|\nabla G(\mathbf{X})\| = \frac{1}{\sqrt{n}}\|\nabla F(\bar{\mathbf{x}})\|$. Therefore, for any $\mathbf{X} \in U$ where $G(\mathbf{X}) \neq G(\mathbf{X}^\infty)$, we can chain these relations and obtain

$$\sqrt{n}\|\nabla G(\mathbf{X})\| \geq \kappa_F |G(\mathbf{X}) - G(\mathbf{X}^\infty)|^{\theta_F}.$$

Rearranging, we obtain the final KL inequality for $G$:

$$\|\nabla G(\mathbf{X})\| \geq \frac{\kappa_F}{\sqrt{n}}|G(\mathbf{X}) - G(\mathbf{X}^\infty)|^{\theta_F}.$$

This indicates $G$ satisfies the KL property at $\mathbf{X}^\infty$ with exponent $\theta_G = \theta_F$ and parameter $\kappa_G = \frac{\kappa_F}{\sqrt{n}}$. $\qquad \square$

Define the following notations:

$$(2.29) \qquad P^t = P(\mathbf{x}^t, \mathbf{v}^t), \quad \Delta P^t = P^t - P^\infty,$$

$$(2.30) \qquad \mathcal{T}^t = \sqrt{\gamma \min\{1, \alpha\}} \left( \|\mathbf{V}^t\|^2 + \|\mathbf{V}^t - \tilde{\mathbf{M}}\mathbf{V}^t\|^2 + \|\mathbf{X}^t - \tilde{\mathbf{M}}\mathbf{X}^t\|^2 \right)^{1/2},$$

where $\gamma$ is the constant from Theorem 2.5. Then, it can be readily verified from (2.26) that

$$(2.31) \qquad\qquad (\mathcal{T}^t)^2 \leq \Delta P^t - \Delta P^{t+1}.$$

To prove the whole sequence convergence of $\{\mathbf{X}^t\}$ generated by UDNA, we will show that the distance sequence $\{\|\mathbf{X}^{t+1} - \mathbf{X}^t\|\}$ is summable, i.e., $\sum_{t=0}^\infty \|\mathbf{X}^{t+1} - \mathbf{X}^t\| \leq \infty$, which means that $\{\mathbf{X}^t\}$ is a Cauchy sequence. The following lemma reveals that $\|\mathbf{X}^{t+1} - \mathbf{X}^t\|$ can be controlled by $\mathcal{T}^t$.

**Lemma 2.10.** *Under the same assumptions and parameter choices of Theorem 2.5, it holds that*

$$(2.32) \qquad\qquad \|\mathbf{X}^{t+1} - \mathbf{X}^t\| \leq \sqrt{c_1}\mathcal{T}^t,$$

*where $c_1 = \frac{2\max\{\alpha^2\Psi, 4\}}{\gamma\min\{1, \alpha\}}$.*

*Proof.* With newly given notations (2.28), it is convenient to rewrite the recursion of $\mathbf{X}^t$ as the following formular,

$$\mathbf{X}^{t+1} = \tilde{\mathbf{A}}\mathbf{X}^t - \alpha\tilde{\mathbf{B}}\mathbf{D}^t,$$

where $\mathbf{D}^t = [(\mathbf{d}_1^t)^\mathsf{T}; (\mathbf{d}_2^t)^\mathsf{T}; \ldots; (\mathbf{d}_n^t)^\mathsf{T}]$ with $\mathbf{d}_i^t = \mathbf{H}_i^t\mathbf{v}_i^t$, $i = 1, \ldots, n$. Note $\mathbf{X}^{t+1} - \mathbf{X}^t = (\tilde{\mathbf{A}} - \mathbf{I}_n)(\mathbf{X}^t - \tilde{\mathbf{M}}\mathbf{X}^t) - \alpha\tilde{\mathbf{B}}\mathbf{D}^t$ and $\rho(\tilde{\mathbf{A}} - \mathbf{I}_n) < 2$. Then it follows from Lemma A.2 and Assumption 4 that

$$\|\mathbf{X}^{t+1} - \mathbf{X}^t\| \leq \sqrt{8\|\mathbf{X}^t - \tilde{\mathbf{M}}\mathbf{X}^t\|^2 + 2\alpha^2\Psi^2\|\mathbf{V}^t\|^2} \leq \sqrt{c_1}\mathcal{T}^t,$$

where $c_1 = \frac{2\max\{\alpha^2\Psi^2, 4\}}{\gamma\min\{1,\alpha\}}$. $\qquad\square$

Therefore, by Lemma 2.10, it suffices to show that the sequence $\{\mathcal{T}^t\}$ is summable. To this end, we first establish a key relationship between $P^{t+1}$ and $G(\mathbf{X}^t)$, and derive an upper bound for $\|\nabla G(\mathbf{X}^t)\|$ in terms of $\mathcal{T}^t$, as summarized in the following lemma.

**Lemma 2.11.** *Under the same assumptions and parameter choices of Theorem 2.5, two inequalities holds as follows,*

$$(2.33) \qquad \Delta P^{t+1} \leq G(\mathbf{X}^t) + c_2(\mathcal{T}^t)^2 - P^\infty$$

*and*

$$(2.34) \qquad \|\nabla G(\mathbf{X}^t)\|^2 \leq (c_3\mathcal{T}^t)^2,$$

*where* $c_2 = \max\left\{1, \frac{(1-\sigma_A^2)(1-\sigma_C^2)}{4(1+\sigma_C^2)L^2\sigma_D^2}\right\} / (\gamma\min\{1,\alpha\})$ *and* $c_3 = \sqrt{\frac{3\max\{1, L^2\}}{\gamma\min\{1,\alpha\}n^2}}$.

*Proof.* Leveraging the descent of $P^t$ and the definition of $\mathcal{T}^t$ (2.30) yields

$$\Delta P^{t+1} \leq P^t - P^\infty$$

$$\leq F((1/n)\mathbf{1}_n^\mathsf{T}\mathbf{X}^t) + \|\mathbf{X}^t - \tilde{\mathbf{M}}\mathbf{X}^t\|^2 + \frac{(1-\sigma_A^2)(1-\sigma_C^2)}{4(1+\sigma_C^2)L^2\sigma_D^2}\|\mathbf{V}^t - \tilde{\mathbf{M}}\mathbf{V}^t\|^2 - P^\infty$$

$$\leq G(\mathbf{X}^t) + c_2(\mathcal{T}^t)^2 - P^\infty,$$

where $c_2 = \max\left\{1, \frac{(1-\sigma_A^2)(1-\sigma_C^2)}{4(1+\sigma_C^2)L^2\sigma_D^2}\right\} / (\gamma\min\{1,\alpha\})$. So the inequality (2.33) is proved. Note the relationship that $\nabla G(\mathbf{X}^t) = (1/n)\mathbf{1}_n\nabla F(\bar{\mathbf{x}}^t)$ and the decomposition that $\mathbf{1}\nabla F(\bar{\mathbf{x}}^t)^\mathsf{T} = \mathbf{1}\nabla F(\bar{\mathbf{x}}^t)^\mathsf{T} - \mathbf{1}\overline{\nabla}f(\mathbf{x}^t)^\mathsf{T} + \mathbf{1}\overline{\nabla}f(\mathbf{x}^t)^\mathsf{T} - \mathbf{V}^t + \mathbf{V}^t$. Then, by Lemma A.2 and Assumption 4, it holds that

$$\|\nabla G(\mathbf{X}^t)\|^2 = \frac{1}{n}\|\nabla F(\bar{\mathbf{x}}^t)\|^2$$

$$\leq \frac{3L^2}{n^2}\|\mathbf{X}^t - \tilde{\mathbf{M}}\mathbf{X}^t\|^2 + \frac{3}{n^2}\|\mathbf{V}^t - \tilde{\mathbf{M}}\mathbf{V}^t\|^2 + \frac{3}{n^2}\|\mathbf{V}^t\|^2 \leq (c_3\mathcal{T}^t)^2,$$

where $c_3 = \sqrt{\frac{3\max\{1, L^2\}}{\gamma\min\{1,\alpha\}n^2}}$. So the inequality (2.34) is proved. $\qquad\square$

Assuming that $F$ satisfies the KL property, Lemma 2.9 guarantees that $G$ inherits this property. From inequality (2.34), we first bound $G(\mathbf{X}^t) - P^\infty$ by $\mathcal{T}^t$. This result, together with the relation given in (2.33), leads directly to an upper bound for $\Delta P^{t+1}$ in terms $\mathcal{T}^t$, as we formally present in the following lemma.

**Lemma 2.12.** *Under the same assumptions and parameter choices of Theorem 2.5, let $\mathbf{X}^\infty = \mathbf{1}_n \otimes (\bar{\mathbf{x}}^\infty)^\mathsf{T}$ be an accumulation point of $\mathbf{X}^t$ where $\bar{\mathbf{x}}^\infty$ is some critical point of $F$. Assume $F$ satisfies the KŁ property at $\bar{\mathbf{x}}^\infty$ with exponent $\theta \in [0, 1)$ and parameters $\{\sqrt{n}\kappa > 0, v = 1\}$. Then, there exists a neighborhood $\mathcal{N}_\infty$ of $\mathbf{X}^\infty$ and some $t_1 \in \mathbb{N}_+$ such that for all $t \in \mathcal{S} := \{t \geq t_1 : \mathbf{X}^t \in \mathcal{N}_\infty\} \neq \emptyset$,*

$$c_3 \mathcal{T}^t < 1, \quad |G(\mathbf{X}^t) - P^\infty| < 1,$$

*and for $\theta \in (0, 1)$ and $G(\mathbf{X}^t) - P^\infty \neq 0$,*

$$(2.35) \qquad \Delta P^{t+1} \leq \kappa^{-\frac{1}{\theta}}(c_3 \mathcal{T}^t)^{\frac{1}{\theta}} + c_2(\mathcal{T}^t)^2,$$

*where $c_2 = \max\left\{1, \frac{(1-\sigma_A^2)(1-\sigma_C^2)}{4(1+\sigma_C^2)L^2\sigma_D^2}\right\}/(\gamma \min\{1, \alpha\})$ and $c_3 = \sqrt{\frac{3\max\{1, L^2\}}{\gamma \min\{1, \alpha\}n^2}}$.*

*Proof.* Since $\mathbf{X}^\infty$ is an accumulation point of $\{\mathbf{X}^t\}$, we have $P^\infty = G(\mathbf{X}^\infty)$. By the KŁ assumption on $F$ and Lemma 2.9, we have $G$ possesses KŁ property at $\mathbf{X}^\infty$ with exponent $\theta \in [0, 1)$ and parameters $\{\kappa > 0, v = 1\}$. Hence, there exists a neighborhood $\mathcal{N}_\infty$ of $\mathbf{X}^\infty$ such that for all $\mathbf{X} \in \mathcal{N}_\infty \cap \{\mathbf{X} : 0 < |G(\mathbf{X}) - P^\infty| < 1\}$,

$$(2.36) \qquad |G(\mathbf{X}) - P^\infty|^\theta \leq \frac{\|\nabla G(\mathbf{X})\|}{\kappa}.$$

From Theorem 2.5, $\mathcal{T}^t \to 0$, $G(\mathbf{X}^t) \to P^\infty$, so there exists some $t_1 \in \mathbb{N}_+$ such that for all $t \geq t_1$

$$c_3 \mathcal{T}^t < 1, \quad |G(\mathbf{X}^t) - P^\infty| < 1,$$

where $c_3 = \sqrt{\frac{3\max\{1, L^2\}}{\gamma \min\{1, \alpha\}n^2}}$. Define $\mathcal{S} = \{t \geq t_1 : \mathbf{X}^t \in \mathcal{N}_\infty\}$ which is nonempty since $\mathbf{X}^\infty$ is an accumulation point of $\{\mathbf{X}^t\}$. For $\theta \in (0, 1)$ and any $t \in \mathcal{S}$ with $G(\mathbf{X}^t) - P^\infty \neq 0$, the inequality (2.36) implies

$$(2.37) \qquad |G(\mathbf{X}^t) - P^\infty| \leq \left(\frac{\|\nabla G(\mathbf{X}^t)\|}{\kappa}\right)^{1/\theta}.$$

Substituting (2.34) into (2.37) yields

$$(2.38) \qquad |G(\mathbf{X}^t) - P^\infty| \leq \kappa^{-\frac{1}{\theta}}(c_3 \mathcal{T}^t)^{\frac{1}{\theta}}.$$

Combining (2.38) with (2.33) gives the desired bound (2.35). □

In the setting of Lemma 2.12, we may assume, without loss of generality, that $|\mathcal{S}| = \infty$. Using (2.31) and Lemma A.3 with $a_1 = \Delta P^t$ and $a_2 = \Delta P^{t+1}$ gives

$$(2.39) \qquad (\mathcal{T}^t)^2 \leq \Delta P^t - \Delta P^{t+1} \leq \frac{1}{1-\theta}(\Delta P^t)^\theta \underbrace{[(\Delta P^t)^{1-\theta} - (\Delta P^{t+1})^{1-\theta}]}_{\Delta P_\theta^t}.$$

Based on the above preliminaries, we now establish the summability of $\{\mathcal{T}^t\}$.

**Theorem 2.13.** *Suppose that Assumptions 1, 2, 3, and 4 hold. Let $\{\mathbf{X}^t\}$ be the sequence generated by Algorithm 1 (UDNA) using the stepsize*

$$\alpha \leq \min\left\{\frac{(1-\sigma_A^2)(1-\sigma_C^2)^2 n}{64L^2\sigma_D^2\Psi^2}, \frac{(1-\sigma_A^2)\psi}{(2L+n+8\sigma_B^2 n)\Psi^2}, \frac{(1-\sigma_A^2)n}{2L^2\Psi^2}\right\}.$$

*Let $\mathbf{X}^\infty = \mathbf{1}_n \otimes (\bar{\mathbf{x}}^\infty)^\mathsf{T}$ be an accumulation point of $\mathbf{X}^t$, where $\bar{\mathbf{x}}^\infty$ is some critical point of $F$. If $F$ satisfies the KŁ property at $\bar{\mathbf{x}}^\infty$ with exponent $\theta \in [0, 1)$ and parameters $\{\sqrt{n}\kappa > 0, v = 1\}$, then $\{\mathcal{T}^t\}$ is summable and the whole sequence $\{\mathbf{X}^t\}$ is convergent.*

*Proof.* By Lemma 2.10, $\|\mathbf{X}^{t+1} - \mathbf{X}^t\| \leq \sqrt{c_1}\mathcal{T}^t$. Thus, summability of $\{\mathcal{T}^t\}$ implies $\{\mathbf{X}^t\}$ is a Cauchy sequence and hence convergent. We prove summability by considering three cases based on the KŁ exponent $\theta$. At first, we define $\mathcal{S}^o = \mathcal{S} \cap \{t : G(\mathbf{X}^t) - P^\infty \neq 0\}$.

**Case I**: $\theta \in (0, 1/2]$. By Lemma 2.12, $c_3\mathcal{T}^t < 1$ for $t \in \mathcal{S}^o$. Since $1/\theta \geq 2$, inequality (2.35) implies

$$(2.40) \qquad \Delta P^{t+1} \leq c_4 (\mathcal{T}^t)^2,$$

where $c_4 = \kappa^{-\frac{1}{\theta}} c_3^2 + c_2$. Computing (2.31) $+ \omega \times$ (2.40), we obtain

$$(1 + \omega)\Delta P^{t+1} \leq \Delta P^t - (1 - \omega c_4)(\mathcal{T}^t)^2.$$

Choosing $\omega \in (0, 1/c_4)$ ensures that

$$(2.41) \qquad \Delta P^{t+1} \leq \frac{1}{1+\omega}\Delta P^t \leq \Delta P^0 \left(\frac{1}{1+\omega}\right)^{t+1}.$$

Applying (2.41) to (2.31) yields

$$(2.42) \qquad \mathcal{T}^t \leq \sqrt{\Delta P^0}(\tau_1)^t,$$

where $\tau_1 = 1/\sqrt{1+\omega}$.

**Case II**: $\theta \in (1/2, 1)$. From (2.35) and $c_3\mathcal{T}^t < 1$ for $t \in \mathcal{S}^o$, we have

$$(2.43) \qquad \Delta P^{t+1} \leq c_5 (\mathcal{T}^t)^{1/\theta},$$

where $c_5 = (c_3/\kappa)^{1/\theta} + c_2$. Using (2.39), (2.43), and Young's inequality with parameter $\eta = 1$, we obtain

$$(2.44)$$

$$\mathcal{T}^{t+1} \leq \sqrt{\frac{1}{1-\theta}(\Delta P^{t+1})^\theta \Delta P_\theta^{t+1}} \leq \sqrt{\frac{c_5^\theta}{1-\theta}\mathcal{T}^t \Delta P_\theta^{t+1}} \leq \frac{1}{2}\mathcal{T}^t + \frac{c_5^\theta}{2 - 2\theta}\Delta P_\theta^{t+1}.$$

**Case III**: $\theta = 0$. For any $t \in \mathcal{S}^o$, the KŁ property of $G$ at $\mathbf{X}^\infty$ impies

$$(2.45) \qquad \kappa(G(\mathbf{X}^t) - P^\infty)^0 \leq \|\nabla G(\mathbf{X}^{t'})\|.$$

Combining with $\|\nabla G(\mathbf{X}^{t'})\|^2 \leq c_3^2 (\mathcal{T}^{t'})^2$ from (2.34) further gives

$$(2.46) \qquad \mathcal{T}^t \geq \frac{\kappa}{c_3}, \ \forall\, t \in \mathcal{S}^o.$$

By (2.31) and the fact that $\Delta P^t \to 0$ monotonically, we must have $\mathcal{T}^t \to 0$. This contradicts the above inequality (2.46) unless $\mathcal{S}^o$ is finite. Thus, there exists some $t_2 \in \mathcal{S}$ such that

$$(2.47) \qquad \mathcal{T}^t = 0, \ \forall\, t \geq t_2,$$

thereby $\mathcal{T}^t$ is summable.

We proceed with analysis following **Case I** and **Case II**. Let $t_4 > t_3$ such that $\tilde{\mathcal{S}} := \{t_3, t_3 + 1, \ldots, t_4 - 1\} \subset \mathcal{S}^o$. Combining (2.42) with (2.44) and summing over $\tilde{\mathcal{S}}$, we obtain

$$(2.48) \qquad \sum_{t=t_3}^{t_4-1} \mathcal{T}^t \leq \max\left\{\frac{\sqrt{\Delta P^0}\tau_1^{t_3}}{1 - \tau_1}, 2\mathcal{T}^{t_3} + \frac{c_5^\theta}{1-\theta}(\Delta P^{t_3})^{1-\theta}\right\},$$

where the second term in R.H.S of the above inequality is obtained by

$$2\sum_{t=t_3}^{t_4-1}\mathcal{T}^{t+1} \le \sum_{t=t_3}^{t_4-1}\mathcal{T}^t + \frac{c_5^\theta}{1-\theta}\sum_{t=t_3}^{t_4-1}((\Delta P^t)^{1-\theta} - (\Delta P^{t+1})^{1-\theta})$$

$$\implies$$

$$\sum_{t=t_3}^{t_4-1}\mathcal{T}^{t+1} \le \mathcal{T}^{t_3} - \mathcal{T}^{t_4} + \frac{c_5^\theta}{1-\theta}(\Delta P^{t_3})^{1-\theta}.$$

Let $r$ be arbitrarily small such that

$$\mathcal{B}_r(\mathbf{X}^\infty) := \{\mathbf{X} : \|\mathbf{X} - \mathbf{X}^\infty\| \le r\} \subseteq \mathcal{N}_\infty.$$

By Theorem 2.5 and the fact that $\mathbf{1}_n \otimes (\bar{\mathbf{x}}^\infty)^\mathsf{T} = \mathbf{X}^\infty$ is an accumulation point of $\mathbf{X}^t$, there exists $t_3' \in \mathcal{S}^o$ such that

$$\|\mathbf{X}^{t_3'} - \mathbf{X}^\infty\| < \frac{r}{2}, \quad \sqrt{c_1}\max\left\{\frac{\sqrt{\Delta P^0}\tau_1^{t_3'}}{1-\tau_1}, 2\mathcal{T}^{t_3'} + \frac{c_5^\theta}{1-\theta}(\Delta P^{t_3'})^{1-\theta}\right\} < \frac{r}{2}.$$

Setting $t_3 = t_3'$ in (2.48) gives $\mathbf{X}^{t_3} \in \mathcal{B}_r(\mathbf{X}^\infty)$ and

(2.49) $$\sum_{t=t_3}^{t_4-1}\mathcal{T}^t \le \frac{r}{2\sqrt{c_1}}.$$

Next, we prove by contradiction that $\mathbf{X}^t \in \mathcal{B}_r(\mathbf{X}^\infty) \subseteq \mathcal{N}_\infty$ for all $t \ge t_3$. Suppose that there exists $t_4' > t_3$ (the smallest such index) such that $\|\mathbf{X}^{t_4'} - \mathbf{X}^\infty\| \ge r$. Then for $t \in \mathcal{S}^o \cap [t_3, t_4')$, setting $t_4 = t_4'$ yields

$$\|\mathbf{X}^{t_4} - \mathbf{X}^\infty\| \le \|\mathbf{X}^{t_3} - \mathbf{X}^\infty\| + \sum_{t=t_3}^{t_4-1}\|\mathbf{X}^{t+1} - \mathbf{X}^t\| < \frac{r}{2} + \sqrt{c_1}\sum_{t=t_3}^{t_4-1}\mathcal{T}^t < r,$$

where the second and last inequalities follow respectively from (2.32) and (2.49). This contradicts the assumption, so $\mathbf{X}^t \in \mathcal{B}_r(\mathbf{X}^\infty)$ for all $t \ge t_3$. Passing the limit $t_4 \to \infty$ to the both sides of (2.49) yields the summability of $\{\mathcal{T}^t\}$, as is $\{\|\mathbf{X}^{t+1} - \mathbf{X}^t\|\}$. Hence, $\{\mathbf{X}^t\}$ is a Cauchy sequence and therefore is convergent. $\square$

A direct consequence of Theorem 2.13 is the following estimations of convergence rates, derived using techniques similar to [6].

**Corollary 2.14.** *Under the setting of Theorem 2.13, the following convergence rates hold:*

1. *When $\theta \in (1/2, 1)$, there exists some $d_1 > 0$ such that*

$$\|\mathbf{X}^t - \mathbf{X}^\infty\| \le d_1 t^{-\frac{1-\theta}{2\theta-1}}, \ \forall \ t \ge 0.$$

2. *When $\theta \in (0, 1/2]$, there exists some $d_2 > 0$ such that*

$$\|\mathbf{X}^t - \mathbf{X}^\infty\| \le d_2(\tau_1)^t, \ \forall \ t \ge 0,$$

*where $\tau_1 = 1/\sqrt{1+\omega}$, $\omega \in (0, 1/c_4)$, and*

$$c_4 = \frac{3\kappa^{-\frac{1}{\theta}}\max\{1, L^2\} + n^2\max\left\{1, \frac{(1-\sigma_A^2)(1-\sigma_C^2)}{4(1+\sigma_C^2)L^2\sigma_D^2}\right\}}{n^2\gamma\min\{1, \alpha\}}.$$

3. *When $\theta = 0$, $\mathbf{X}^t$ coverges to $\mathbf{X}^\infty$ in a finite number of steps.*

2.4. **Methods of constructing Hessian inverse approximations.** We begin with reviewing how previous decentralized quasi-Newton methods ensured bounded approximate matrices for the Hessian or its inverse. It has been observed that the Hessian inverse approximations generated by (2.7) in DQN [47] may exhibit negative eigenvalues, despite being theoretically assumed to maintain positive eigenvalues with lower boundedness. To establish convergence guarantees for DGM-BB-C [18] and DSG [25] in nonconvex optimization settings, it becomes necessary to project the BB-type matrix onto the space of positive definite matrices. While the regularization and damping techniques proposed in [62] for ensuring eigenvalue boundedness were initially employed for convex optimization, we emphasize that these methods are fundamentally independent of convexity assumptions and therefore directly applicable to nonconvex scenarios. Below, we propose several novel approaches for constructing bounded approximate matrices of Hessian inverse, specifically designed for nonconvex optimization frameworks.

2.4.1. *Memoryless quasi-Newton updating schemes.* Empirical evidence from comparisons between memoryless quasi-Newton methods and standard BFGS [5] indicates that the accuracy of the Hessian approximation is not essential for their efficiency and robustness. Inspired by this finding and the need for low computational cost, we study two specific memoryless quasi-Newton updating schemes in decentralized optimization.

*Memoryless SR1.* The memoryless SR1 method with direct approximation to the Hessian is obtained by considering $\mathbf{H}_i^t = \mathbf{I}_p$ in the standard SR1 update, i.e.,

$$(2.50) \qquad \hat{\mathbf{H}}_i^{t+1} = \mathbf{I}_p + \frac{(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)^\mathsf{T}}{(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t},$$

where $\mathbf{s}_i^t = \mathbf{x}_i^{t+1} - \mathbf{x}_i^t$ and $\check{\mathbf{y}}_i^t = \mathbf{v}_i^{t+1} - \mathbf{v}_i^t$. The main drawback of memoryless SR1 update (2.50) is that the SR1 approximation may not be positive definite along the iterations. Thus, a necessary safeguarding on (2.50) is proposed:

$$(2.51) \qquad \mathbf{H}_i^{t+1} = \begin{cases} \mathbf{I}_p + \frac{(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)^\mathsf{T}}{(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t}, & \text{if } [\lambda_{\min}(\hat{\mathbf{H}}_i^{t+1}), \lambda_{\max}(\hat{\mathbf{H}}_i^{t+1})] \subset [\bar{l}, \bar{u}] \\ & \quad \text{and } (\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t \neq 0; \\ \mathbf{I}_p, & \text{otherwise,} \end{cases}$$

where $0 < \bar{l} \ll \bar{u} < +\infty$ are two constants. As a rank-one perturbation on the identity matrix , the eigenvalue of $\hat{\mathbf{H}}_i^{t+1}$ other than one is easily computed, that is $1 + \frac{\|\mathbf{s}_i^t - \check{\mathbf{y}}_i^t\|^2}{(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t}$. Thus, we easily check if this eigenvalue lies in the interval and compute the bounds on the eigenvalues of $\mathbf{H}_i^{t+1}$. From (2.51), after some simple algebraic manipulations, the memoryless SR1 search direction $\mathbf{d}_i^{t+1} = -\mathbf{H}_i^{t+1}\mathbf{v}_i^{t+1}$, when $[\lambda_{\min}(\hat{\mathbf{H}}_i^{t+1}), \lambda_{\max}(\hat{\mathbf{H}}_i^{t+1})] \subset [\bar{l}, \bar{u}]$ and $(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t \neq 0$, can be written as

$$(2.52) \qquad \mathbf{d}_i^{t+1} = -\mathbf{v}_i^{t+1} - \frac{(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)^\mathsf{T}\mathbf{v}_i^{t+1}}{(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t}(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t).$$

The main advantage of the memoryless SR1 update (2.52) is that for its implementation in computer programs, only three vector-to-vector products, $\|\mathbf{s}_i^t - \check{\mathbf{y}}_i^t\|^2$, $(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)^\mathsf{T}\mathbf{v}_i^{t+1}$ and $(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t$, should be computed. This is very advantageous for solving large-scale problems.

*Memoryless BFGS.* We now let $\Omega_i^t = \left\{ \mathbf{y} \in \mathbb{R}^p : (\mathbf{s}_i^t)^\mathsf{T}\mathbf{y} > 0 \right\}$ and define the function $H_i^t : \Omega_i^t \to \mathbb{R}^{p \times p}$ as

$$(2.53) \qquad H_i^t(\mathbf{y}) = \tau_i^t \left( \mathbf{I}_p - \frac{\mathbf{s}_i^t(\mathbf{y})^\mathsf{T} + \mathbf{y}(\mathbf{s}_i^t)^\mathsf{T}}{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}} \right) + \left( 1 + \frac{\tau_i^t \|\mathbf{y}\|^2}{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}} \right) \frac{\mathbf{s}_i^t(\mathbf{s}_i^t)^\mathsf{T}}{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}},$$

where $\mathbf{s}_i^t = \mathbf{x}_i^{t+1} - \mathbf{x}_i^t$. We propose a memoryless BFGS method to generate $\mathbf{H}^t$ in (2.2),

$$(2.54) \quad \mathbf{H}_i^{t+1} = H_i^t(\mathbf{y}_i^t) = \tau_i^t \left( \mathbf{I}_p - \frac{\mathbf{s}_i^t(\mathbf{y}_i^t)^\mathsf{T} + \mathbf{y}_i^t(\mathbf{s}_i^t)^\mathsf{T}}{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t} \right) + \left( 1 + \frac{\tau_i^t \|\mathbf{y}_i^t\|^2}{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t} \right) \frac{\mathbf{s}_i^t(\mathbf{s}_i^t)^\mathsf{T}}{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t}$$

$$= \frac{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t}{\|\mathbf{y}_i^t\|^2}\mathbf{I}_p - \frac{\mathbf{s}_i^t(\mathbf{y}_i^t)^\mathsf{T} + \mathbf{y}_i^t(\mathbf{s}_i^t)^\mathsf{T}}{\|\mathbf{y}_i^t\|^2} + 2\frac{\mathbf{s}_i^t(\mathbf{s}_i^t)^\mathsf{T}}{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t}.$$

where $\tau_i^t = \frac{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t}{\|\mathbf{y}_i^t\|^2}$,

$$(2.55) \quad \mathbf{y}_i^t = \begin{cases} \check{\mathbf{y}}_i^t, & \text{if } [\lambda_{\min}(H_i^t(\check{\mathbf{y}}_i^t)), \lambda_{\max}(H_i^t(\check{\mathbf{y}}_i^t))] \subset [l, u] \text{ and } (\mathbf{s}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t \neq 0; \\ \hat{\mathbf{y}}_i^t, & \text{otherwise,} \end{cases}$$

$\check{\mathbf{y}}_i^t = \mathbf{v}_i^{t+1} - \mathbf{v}_i^t$, $\hat{\mathbf{y}}_i^t = \mathbf{g}_i^{t+1} - \mathbf{g}_i^t + h_i^t\mathbf{s}_i^t$, $h_i^t = \varrho + \max\left\{ -\frac{(\mathbf{s}_i^t)^\mathsf{T}(\mathbf{g}_i^{t+1} - \mathbf{g}_i^t)}{\|\mathbf{s}_i^t\|^2}, 0 \right\}$, and $\varrho > 0$ is a small constant.

The adaptive selection mechanism in (2.55) is to ensure the quasi-Newton matrix (2.54) has positive bounded eigenvalues. Since $\mathbf{v}_i^t$ is generated by the gradient tracking technique in (2.3), it captures some information of the average gradients on different nodes. So, we prefer using $\mathbf{v}_i^t$ to generate the quasi-Newton matrix by $H_i^t(\check{\mathbf{y}}_i^t)$ whenever possible, where $\check{\mathbf{y}}_i^t = \mathbf{v}_i^{t+1} - \mathbf{v}_i^t$. However, $H_i^t(\check{\mathbf{y}}_i^t)$ is not necessarily positive definite and bounded. So, in (2.55) we simply check if the smallest and largest eigenvalues of the quasi-Newton matrix $H_i^t(\check{\mathbf{y}}_i^t)$ belong to an interval $[l, u]$, where $0 < l \ll u$ are two parameters. If not, we would use the alternative $\mathbf{g}_i^t$ with correction $h_i^t\mathbf{s}_i^t$ to update the quasi-Newton matrix by $H_i^t(\hat{\mathbf{y}}_i^t)$, where $\hat{\mathbf{y}}_i^t = \mathbf{g}_i^{t+1} - \mathbf{g}_i^t + h_i^t\mathbf{s}_i^t$. By the BFGS updating formula (2.54), $H_i^t(\hat{\mathbf{y}}_i^t)$ is guaranteed to be positive definite and uniformly bounded. We now explain the computation cost of the smallest and largest eigenvalues of $\mathbf{H}_i^{t+1}$ is in fact negligible. Note that $\mathbf{H}_i^{t+1} = H_i^t(\check{\mathbf{y}}_i^t)$ is obtained from the scalar matrix $\tau_i^t\mathbf{I}_p$ by a rank-two BFGS update. Hence, if $\tau_i^t = \frac{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t}{\|\mathbf{y}_i^t\|^2} > 0$, $\mathbf{H}_i^{t+1}$ has $p-2$ eigenvalues of $\tau_i^t$ and two eigenvalues $0 < \lambda_i^{t+1} \leq \Lambda_i^{t+1}$ satisfying

$$(2.56) \qquad \begin{cases} \lambda_i^{t+1}\Lambda_i^{t+1} & = \frac{\|\mathbf{s}_i^t\|^2}{\|\mathbf{y}_i^t\|^2}, \\ \lambda_i^{t+1} + \Lambda_i^{t+1} & = \frac{2\|\mathbf{s}_i^t\|^2}{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t}. \end{cases}$$

The system (2.56) defines a quadratic equation, which has two roots

$$(2.57) \qquad \begin{cases} \lambda_i^{t+1} = \lambda_{\min}(H_i^t(\mathbf{y}_i^t)) = \frac{\|\mathbf{s}_i^t\|^2}{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t} \left( 1 - \sqrt{1 - \frac{((\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t)^2}{\|\mathbf{s}_i^t\|^2\|\mathbf{y}_i^t\|^2}} \right), \\ \Lambda_i^{t+1} = \lambda_{\max}(H_i^t(\mathbf{y}_i^t)) = \frac{\|\mathbf{s}_i^t\|^2}{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t} \left( 1 + \sqrt{1 - \frac{((\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t)^2}{\|\mathbf{s}_i^t\|^2\|\mathbf{y}_i^t\|^2}} \right). \end{cases}$$

$\mathbf{H}^t$ given as (2.54) is easily shown to be bounded with $\tau_i^t = \frac{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t}{\|\mathbf{y}_i^t\|^2}$. From (2.55), we acquire $l\mathbf{I}_p \preceq \mathbf{H}_i^{t+1} \preceq u\mathbf{I}_p$ when $\mathbf{y}_i^t$ takes $\check{\mathbf{y}}_i^t$. When $\mathbf{y}_i^t$ takes $\hat{\mathbf{y}}_i^t$, we have by

(2.57) that

$$\frac{(\mathbf{s}_i^t)^\mathsf{T}\hat{\mathbf{y}}_i^t}{2\|\hat{\mathbf{y}}_i^t\|^2}\mathbf{I}_p \preceq \mathbf{H}_i^{t+1} \preceq \frac{2\|\mathbf{s}_i^t\|^2}{(\mathbf{s}_i^t)^\mathsf{T}\hat{\mathbf{y}}_i^t}\mathbf{I}_p.$$

According to $\|\hat{\mathbf{y}}_i^t\| = \|\mathbf{g}_i^{t+1} - \mathbf{g}_i^t + h_i^t\mathbf{s}_i^t\| \leq (2L + \varrho)\|\mathbf{s}_i^t\|$ yielded by the $L$-Lipschitz continuity of $\nabla f_i$, we have $(\mathbf{s}_i^t)^\mathsf{T}\hat{\mathbf{y}}_i^t \geq \frac{1}{L}\|\hat{\mathbf{y}}_i^t\|^2 + \varrho\|\mathbf{s}_i^t\|^2 \geq \hat{\varrho}\|\hat{\mathbf{y}}_i^t\|^2$ with $\hat{\varrho} = \frac{1}{4L^2+4L\varrho+\varrho^2} + \frac{1}{L}$, which further implies the left hand side of the above inequality can be bounded below by $\frac{\hat{\varrho}}{2}$. By $(\mathbf{s}_i^t)^\mathsf{T}\hat{\mathbf{y}}_i^t \geq \varrho\|\mathbf{s}_i^t\|^2$, the right hand side of the above inequality can be bounded above by $\frac{2}{\varrho}$. Hence, we obtain

$$\min\left\{l, \frac{\hat{\varrho}}{2}\right\}\mathbf{I} \preceq \mathbf{H}^{t+1} \preceq \max\left\{u, \frac{2}{\varrho}\right\}\mathbf{I}.$$

The corresponding search direction of the memoryless BFGS method is $\mathbf{d}_i^{t+1} = -\mathbf{H}_i^{t+1}\mathbf{v}_i^{t+1}$, where $\mathbf{H}_i^{t+1}$ is given by (2.54), i.e.,

$$(2.58) \quad \mathbf{d}_i^{t+1} = -\frac{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t}{\|\mathbf{y}_i^t\|^2}\mathbf{v}_i^{t+1} + \frac{((\mathbf{y}_i^t)^\mathsf{T}\mathbf{v}_i^{t+1})\mathbf{s}_i^t + ((\mathbf{s}_i^t)^\mathsf{T}\mathbf{v}_i^{t+1})\mathbf{y}_i^t}{\|\mathbf{y}_i^t\|^2} - 2\frac{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{v}_i^{t+1}}{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t}\mathbf{s}_i^t.$$

Observe that the numerical computation of $\mathbf{d}_i^{t+1}$ from (2.58) involves only five scalar products: $\|\mathbf{s}_i^t\|^2$, $\|\mathbf{y}_i^t\|^2$, $(\mathbf{s}_i^t)^\mathsf{T}\mathbf{y}_i^t$, $(\mathbf{s}_i^t)^\mathsf{T}\mathbf{v}_i^{t+1}$, and $(\mathbf{y}_i^t)^\mathsf{T}\mathbf{v}_i^{t+1}$. Therefore, it is very suitable for applying the memoryless BFGS method to solve large-scale problems.

2.4.2. *Corrected quasi-Newton updating schemes.* Motivated by [11], we introduce a useful correcting method for the variation of average-gradient approximation $\check{\mathbf{y}}_i^t = \mathbf{v}_i^{t+1} - \mathbf{v}_i^t$,

$$(2.59) \qquad\qquad \check{\mathbf{y}}_i^t = \eta_i^t\check{\mathbf{y}}_i^t + (1 - \eta_i^t)\mathbf{s}_i^t,$$

where $\eta_i^t \in (0, 1]$ such that

$$(2.60) \qquad\qquad \frac{(\mathbf{s}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t}{\|\mathbf{s}_i^t\|^2} \geq \lambda$$

for some $\lambda \in (0, 1)$. The positivity condition (2.60), which is crucial for ensuring positive definite Hessian inverse approximations in both centralized and decentralized settings, can be ensured by choosing a sufficiently small $\eta_i^t \in (0, 1]$ in (2.59). In the following, we provide an adaptive approach for generating $\eta_i^t$,

$$(2.61) \qquad\qquad \eta_i^t = \min\left\{\hat{\eta}_i^t, \frac{\hat{L}\|\mathbf{s}_i^t\|}{\|\check{\mathbf{y}}_i^t\|}\right\}$$

where $\hat{L}$ is some positive constant and $\hat{\eta}_i^t$ is defined as

$$(2.62) \qquad\qquad \hat{\eta}_i^t = \begin{cases} \frac{(1-\lambda)\|\mathbf{s}_i^t\|^2}{\|\mathbf{s}_i^t\|^2 - (\mathbf{s}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t}, & \text{if } (\mathbf{s}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t \leq \lambda\|\mathbf{s}_i^t\|^2; \\ 1, & \text{otherwise.} \end{cases}$$

Therefore, the positivity of $(\mathbf{s}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t$ in (2.60) can be guaranteed as the following lemma shows.

**Lemma 2.15.** *For $\eta_i^t$ defined in (2.61) and $\check{\mathbf{y}}_i^t$ defined in (2.59), we have $0 < \eta^t \leq 1$ and $(\mathbf{s}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t \geq \lambda\|\mathbf{s}_i^t\|^2$.*

*Proof.* When $(\mathbf{s}_i^t)^{\mathsf{T}}\breve{\mathbf{y}}_i^t \le \lambda\|\mathbf{s}_i^t\|^2$, it can be readily verified that $0 < \hat{\eta}_i^t \le 1$. Since $\eta_i^t = \min\{\hat{\eta}_i^t, \frac{\hat{L}\|\mathbf{s}_i^t\|}{\|\breve{\mathbf{y}}_i^t\|}\}$, we obtain $0 < \eta^t \le 1$. Multiplying left and right sides of (2.59) by $(\mathbf{s}_i^t)^{\mathsf{T}}$ and using the expression of $\eta_i^t$ from (2.61), we have

$$(\mathbf{s}_i^t)^{\mathsf{T}}\breve{\mathbf{y}}_i^t = \eta_i^t((\mathbf{s}_i^t)^{\mathsf{T}}\breve{\mathbf{y}}_i^t - \|\mathbf{s}_i^t\|^2) + \|\mathbf{s}_i^t\|^2 = \begin{cases} \lambda\|\mathbf{s}_i^t\|^2, & \text{if } (\mathbf{s}_i^t)^{\mathsf{T}}\breve{\mathbf{y}}_i^t \le \lambda\|\mathbf{s}_i^t\|^2; \\ ((\mathbf{s}_i^t)^{\mathsf{T}}\breve{\mathbf{y}}_i^t, & \text{otherwise,} \end{cases}$$

which implies that $(\mathbf{s}_i^t)^{\mathsf{T}}\breve{\mathbf{y}}_i^t \ge \lambda\|\mathbf{s}_i^t\|^2$.                                     $\square$

The Dai-Kou [12] and Hager-Zhang [19] conjugate gradient methods are closely related to quasi-Newton methods, as they effectively capture second-order information while maintaining computational efficiency. In fact, these methods can be formally recast as specific instances of quasi-Newton methods. Extending this idea to decentralized optimization, we consider a direction $\tilde{\mathbf{d}}_i^t$ combining the tracking variable $\mathbf{v}_i^t$ with iterate variation $\mathbf{s}_i^{t-1} = \mathbf{x}_i^t - \mathbf{x}_i^{t-1}$:

$$(2.63) \qquad\qquad \tilde{\mathbf{d}}_i^t = -\mathbf{v}_i^t + \beta_i^{t,\tau}\mathbf{s}_i^{t-1},$$

where

$$(2.64) \qquad\qquad \beta_i^{t,\tau} = \frac{(\mathbf{v}_i^t)^{\mathsf{T}}\breve{\mathbf{y}}_i^{t-1}}{(\mathbf{s}_i^{t-1})^{\mathsf{T}}\breve{\mathbf{y}}_i^{t-1}} - \tau\frac{\|\breve{\mathbf{y}}_i^{t-1}\|^2(\mathbf{s}_i^{t-1})^{\mathsf{T}}\mathbf{v}_i^t}{((\mathbf{s}_i^{t-1})^{\mathsf{T}}\breve{\mathbf{y}}_i^{t-1})^2},$$

and the parameter $\tau$ is suggested to lie in the interval $[1, 2]$ according to [19, 12].

*Remark* 2.16. Consider the case with $n = 1$, implying $\tilde{\mathbf{W}} = [1]$, and assume $\eta_i^t \equiv 1$. Then the parameter $\beta_i^{t,\tau}$ in (2.64) is equivalent to Dai-Kou and Hager-Zhang conjugate parameters for $\tau = 1$ and 2, respectively.

Notably, the direction (2.63) can be rewritten as

$$\tilde{\mathbf{d}}_i^t = -\tilde{\mathbf{H}}_i^t\mathbf{v}_i^t,$$

where

$$\tilde{\mathbf{H}}_i^t = \mathbf{I} - \frac{\mathbf{s}_i^{t-1}(\mathbf{z}_i^{t-1})^{\mathsf{T}}}{(\mathbf{s}_i^{t-1})^{\mathsf{T}}\breve{\mathbf{y}}_i^{t-1}}, \quad \mathbf{z}_i^{t-1} = \breve{\mathbf{y}}_i^{t-1} - \tau p_i^{t-1}\mathbf{s}_i^{t-1}, \quad p_i^{t-1} = \frac{\|\breve{\mathbf{y}}_i^{t-1}\|^2}{(\mathbf{s}_i^{t-1})^{\mathsf{T}}\breve{\mathbf{y}}_i^{t-1}}.$$

By symmetrizing $\tilde{\mathbf{H}}_i^t$, we acquire a corrected quasi-Newton update for $\mathbf{H}_i^t$:

$$(2.65) \qquad \mathbf{H}_i^t = \frac{\tilde{\mathbf{H}}_i^t + (\tilde{\mathbf{H}}_i^t)^{\mathsf{T}}}{2} = \mathbf{I} - \frac{1}{2}\frac{\mathbf{s}_i^{t-1}(\mathbf{z}_i^{t-1})^{\mathsf{T}} + \mathbf{z}_i^{t-1}(\mathbf{s}_i^{t-1})^{\mathsf{T}}}{(\mathbf{s}_i^{t-1})^{\mathsf{T}}\breve{\mathbf{y}}_i^{t-1}},$$

termed the Dai-Kou type update for $\tau = 1$ and the Hager-Zhang type update for $\tau = 2$. Sequentially, the search direction $\mathbf{d}_i^t$ for this corrected quasi-Newton method is calculated by

$$\mathbf{d}_i^t = -\mathbf{H}_i^t\mathbf{v}_i^t = -\mathbf{v}_i^t + \frac{(\mathbf{z}_i^{t-1})^{\mathsf{T}}\mathbf{v}_i^t}{2(\mathbf{s}_i^{t-1})^{\mathsf{T}}\breve{\mathbf{y}}_i^{t-1}}\mathbf{s}_i^{t-1} + \frac{(\mathbf{s}_i^{t-1})^{\mathsf{T}}\mathbf{v}_i^t}{2(\mathbf{s}_i^{t-1})^{\mathsf{T}}\breve{\mathbf{y}}_i^{t-1}}\mathbf{z}_i^{t-1}.$$

Now we establish uniform bounds on the eigenvalues of $\mathbf{H}_i^t$ for all $i$ and $t$.

**Lemma 2.17.** *For $\mathbf{H}_i^t$ defined by (2.65), we have*

$$(2.66) \qquad\qquad \lambda_{\max}(\mathbf{H}_i^t) \le \frac{2\tau(\hat{L}^2 + 1)}{\lambda^2},$$

$$(2.67) \qquad\qquad \lambda_{\min}(\mathbf{H}_i^t) \ge \frac{1}{2}.$$

*Proof.* Firstly, it naturally holds $\lambda_{\max}(\mathbf{H}_i^t) \leq \|\mathbf{H}_i^t\|_2$. From Lemma 2.15, we have

$$\|\mathbf{H}_i^t\|_2 = \frac{\tau\|\mathbf{s}_i^{t-1}\|^2\|\check{\mathbf{y}}_i^{t-1}\|^2}{\left((\mathbf{s}_i^{t-1})^{\mathsf{T}}\check{\mathbf{y}}_i^{t-1}\right)^2} \leq \frac{\tau\|\check{\mathbf{y}}_i^{t-1}\|^2}{\lambda(\mathbf{s}_i^{t-1})^{\mathsf{T}}\check{\mathbf{y}}_i^{t-1}} \leq \frac{2\tau(\eta_i^t)^2\|\check{\mathbf{y}}_i^{t-1}\|^2 + 2\tau(1-\eta_i^t)^2\|\mathbf{s}_i^{t-1}\|^2}{\lambda(\mathbf{s}_i^{t-1})^{\mathsf{T}}\check{\mathbf{y}}_i^{t-1}}$$

$$\leq \frac{2\tau\hat{L}^2\|\mathbf{s}_i^{t-1}\|^2 + 2\tau(1-\eta_i^t)^2\|\mathbf{s}_i^{t-1}\|^2}{\lambda(\mathbf{s}_i^{t-1})^{\mathsf{T}}\check{\mathbf{y}}_i^{t-1}} \leq \frac{2\tau(\hat{L}^2+1)}{\lambda^2},$$

where the second inequality applies (2.59) and Lemma A.2, and the third inequality uses the expression of $\eta_i^t$ from (2.61). So the upper bound (2.66) is established.

For any vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^p$, notice that

$$\left(\mathbf{u}\mathbf{v}^{\mathsf{T}} + \mathbf{v}\mathbf{u}^{\mathsf{T}}\right)\left(\mathbf{u} \pm \frac{\|\mathbf{u}\|}{\|\mathbf{v}\|}\mathbf{v}\right) = \left(\mathbf{u}^{\mathsf{T}}\mathbf{v} \pm \|\mathbf{u}\|\|\mathbf{v}\|\right)\left(\mathbf{u} \pm \frac{\|\mathbf{u}\|}{\|\mathbf{v}\|}\mathbf{v}\right).$$

By this, it is not difficult to see that the minimal eigenvalue of $\mathbf{H}_i^t$ is

$$(2.68) \qquad \lambda_{\min}(\mathbf{H}_i^t) = \min\left\{1, 1 - \frac{1}{2}\left(\frac{(\mathbf{s}_i^{t-1})^{\mathsf{T}}\mathbf{z}_i^{t-1}}{(\mathbf{s}_i^{t-1})^{\mathsf{T}}\check{\mathbf{y}}_i^{t-1}} + \frac{\|\mathbf{s}_i^{t-1}\|\|\mathbf{z}_i^{t-1}\|}{|(\mathbf{s}_i^{t-1})^{\mathsf{T}}\check{\mathbf{y}}_i^{t-1}|}\right)\right\}.$$

We define

$$q_i^{t-1} = \frac{\|\check{\mathbf{y}}_i^{t-1}\|^2\|\mathbf{s}_i^{t-1}\|^2}{((\mathbf{s}_i^{t-1})^{\mathsf{T}}\check{\mathbf{y}}_i^{t-1})^2}.$$

By this definition of $q_i^{t-1}$, we can rewrite (2.68) as

$$(2.69) \qquad \lambda_{\min}(\mathbf{H}_i^t) = \frac{1}{2}\left(1 + \tau q_i^{t-1} - \sqrt{\tau^2(q_i^{t-1})^2 - 2\tau q_i^{t-1} + q_i^{t-1}}\right).$$

For the second term in the braces of (2.69), it is monotonically decreasing for $q_i^t \geq 1$ and hence is always greater than its limit $\frac{1}{2}$ as $q_i^t$ tends to $+\infty$. Thus we always have $\lambda_{\min}(\mathbf{H}_i^t) \geq \frac{1}{2}$, which proves the lower bound (2.67). $\square$

## 3. NUMERICAL EXPERIMENTS

In this section, we would like to examine the performance of our developed algorithms as the following outline:

a. Compare UDNAs of different approximations $\mathbf{H}^t$s to Hessian inverse;
b. Compare UDNAs with some well-developed nonconvex optimization algorithms, including gradient-based algorithms [36, 1, 50, 20] and quasi-Newton algorithms [47, 62].

The considered optimization problem is smooth but nonconvex, over a connected undirected network with edge density $d \in (0, 1]$. For the generated network, we choose the Metropolis constant edge weight matrix [54] as the mixing matrix, that is

$$\tilde{W}_{ij} = \begin{cases} \frac{1}{\max\{\deg(i), \deg(j)\}+1}, & \text{if } (i, j) \in \mathcal{E}; \\ 0, & \text{if } (i, j) \notin \mathcal{E} \text{ and } i \neq j; \\ 1 - \sum_{k \in \mathcal{N}_i/\{i\}} \tilde{W}_{ik}, & \text{if } i = j, \end{cases}$$

where $(i, j) \in \mathcal{E}$ indicates there is an edge between node $i$ and node $j$, and $\deg(i)$ means the degree of node $i$. In our experiments, we introduce the communication

volume which can be calculated as follows:

$$\text{Communication volume} = \text{iteration number} \times \text{ the number of edges } \frac{dn(n-1)}{2}$$
$$\times \text{ number of communication rounds per iteration}$$
$$\times \text{ dimension of transmitted vectors on each edge.}$$

In all experiments, we set the number of nodes $n = 10$ and the edge density $d = 0.56$ for the network. For all comparison algorithms, we initialize $\mathbf{x}^0 = \mathbf{0}$. All experiments are coded in MATLAB R2017b and run on a laptop with Intel Core i5-9300H CPU, 16GB RAM, and Windows 10 operating system.

We consider the nonconvex decentralized binary classification problem. Using a logistic regression formulation with a nonconvex regularization, the optimization is given by

$$(3.1) \qquad \min_{\mathbf{z} \in \mathbb{R}^p} \sum_{i=1}^{n} \sum_{j=1}^{n_i} \log\left(1 + \exp(-b_{ij}\mathbf{a}_{ij}^{\mathsf{T}}\mathbf{z})\right) + \hat{\lambda} \sum_{k=1}^{p} \frac{\mathbf{z}_{[k]}^2}{1 + \mathbf{z}_{[k]}^2},$$

where $\mathbf{a}_{ij} \in \mathbb{R}^p$ is the feature vector, $b_{ij} \in \{-1, +1\}$ is the label, $\mathbf{z}_{[k]}$ denotes the $k$-th component of the vector $\mathbf{z}$, and $\hat{\lambda} > 0$ is the regularization parameter. The logistic loss function are semi-algebraic functions [28], and the regularization term as a rational function, is also semi-algebraic. The sum of semi-algebraic functions remains semi-algebraic. All semi-algebraic functions satisfy the KŁ property. Thus, this objective function in (3.1) satisfies the KŁ property on $\mathbb{R}^p$. The experiments are conducted on four datasets in Table 3 from the LIBSVM library: **mushrooms**, **ijcnn1**, **w8a**, **a9a**, **colon-cancer**, and **duke breast-cancer**. The regularization parameter $\hat{\lambda} = 1$.

From the first-order stationarity given in (2.9), the success of each algorithm is measured by the optimality error stated as

$$\left\| \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(\mathbf{x}_i^t) \right\| + \|\mathbf{x}^t - \mathbf{M}\mathbf{x}^t\|.$$

TABLE 3. Datasets

| Dataset | # of samples $\left(\sum_{i=1}^{n} n_i\right)$ | # of features $(p)$ |
|---|---|---|
| **mushrooms** | 8124 | 112 |
| **ijcnn1** | 49990 | 22 |
| **w8a** | 49749 | 300 |
| **a9a** | 32561 | 123 |
| **colon-cancer** | 62 | 2000 |
| **duke breast-cancer** | 44 | 7129 |

We first investigate the performance of UDNAs with different approximations to Hessian inverse, namely different $\mathbf{H}_i^t$s. Based on the proposed techniques in Section 2.3, four variants of UDNA are obtained and represented as UDNA($j$), $j = 1, 2, 3, 4$.

- UDNA(1) implements the memoryless SR1 update. Its detailed procedure is given in Algorithm 2 (Appendix B).

- UDNA(2) corresponds to the memoryless BFGS method, outlined in Algorithm 3 (Appendix B).
- UDNA(3) and UDNA(4) represent Dai-Kou type and Hager-Zhang type corrected quasi-Newton methods, respectively. Both procedures are summarized in Algorithm 4 (Appendix B).

For datasets **mushrooms(ijcnn1;w8a;a9a;colon-cancer;duke breast-cancer)**, algorithm parameters are set as follows their better performance. We set

- $\alpha = 0.12(0.24; 0.21; 0.16; 0.038; 0.022)$, $\bar{l} = 10^{-6}$, and $\bar{u} = 10^6$ in UDNA(1);
- $\alpha = 0.22(0.32; 0.26; 0.34; 0.14; 0.2)$, $\varrho = 0.05(0.05; 0.01; 0.01; 0.3; 0.5)$, $l = 10^{-6}$, and $u = 10^6$ in UDNA(2);
- $\alpha = 0.09(0.14; 0.11; 0.11; 0.013; 0.0072)$, $\lambda = 0.7(0.8; 0.8; 0.8; 0.7; 0.7)$ and $\hat{L} = 1(1; 3; 2; 4; 4)$ in UDNA(3);
- $\alpha = 0.05(0.09; 0.06; 0.07; 0.009; 0.0046)$, $\lambda = 0.7(0.8; 0.8; 1.1; 0.8; 0.8)$ and $\hat{L} = 2(1; 1; 2; 3; 3)$ in UDNA(4).

Fig. 1 shows the decreasing curves of optimality error versus communication for UDNAs on the different datasets. The performance of these algorithms is significantly different. UNDA(2) converges fastest on the datasets with a large number of samples while UNDA(3) wins on datasets with high feature dimensionality. Therefore, we summarize the rankings and average performance of UDNAs across the six datasets in Table 4. UNDA(2) and UNDA(1) respectively performs best and worst. Hence, we take UDNA(1) and UDNA(2) as default of UDNA in the following numerical comparisons with other algorithms.

TABLE 4. Rankings of UDNAs on different datasets

|  | UDNA(1) | UDNA(2) | UDNA(3) | UDNA(4) |
|---|---|---|---|---|
| **mushrooms** | 3 | 2 | 4 | 1 |
| **ijcnn1** | 2 | 1 | 4 | 3 |
| **w8a** | 2 | 1 | 4 | 3 |
| **a9a** | 4 | 1 | 3 | 2 |
| **colon-cancer** | 4 | 3 | 1 | 2 |
| **duke breast-cancer** | 4 | 3 | 1 | 2 |
| average | 3.17 | 1.83 | 2.83 | 2.17 |

Next, we report the numerical result of UDNAs compared with several advanced nonconvex algorithms. Comparison algorithms are listed: Gradient Tracking (GT) [36], Global Update Tracking (GUT) [1], Momentum Tracking (MT) [50], Distributed Stochastic Momentum Tracking (DSMT) [20], UDNA(1), and UDNA(2). Although GUT, MT, and DSMT are stochastic methods, since we are focusing on comparing deterministic decentralized methods, full gradient is used for these methods. For datasets **mushrooms(ijcnn1;w8a;a9a)**, algorithm parameters are set following their better performance and parameter notations follow the source papers. We set $\eta = 0.06(0.09; 0.09; 0.08)$ in GT; set $\eta_t = (0.01; 0.01; 0.01; 0.02) \times n^{1/2}/t^{1/3}$ and $\mu = 0.3$ in GUT; set $\eta = 0.05$ and $\beta = 0.31(0.41; 0.33; 0.35)$ in MT. We set $\eta_w = 1/(1 + \sqrt{1 - (1 - \sigma)^2})$, $\beta = 1 - (1 - \sqrt{\eta_w})/n^{1/3}$, and $\alpha = 0.04(0.08; 0.04; 0.08)$ in DSMT, where $\sigma$ is given by **Lemma 2.2**.
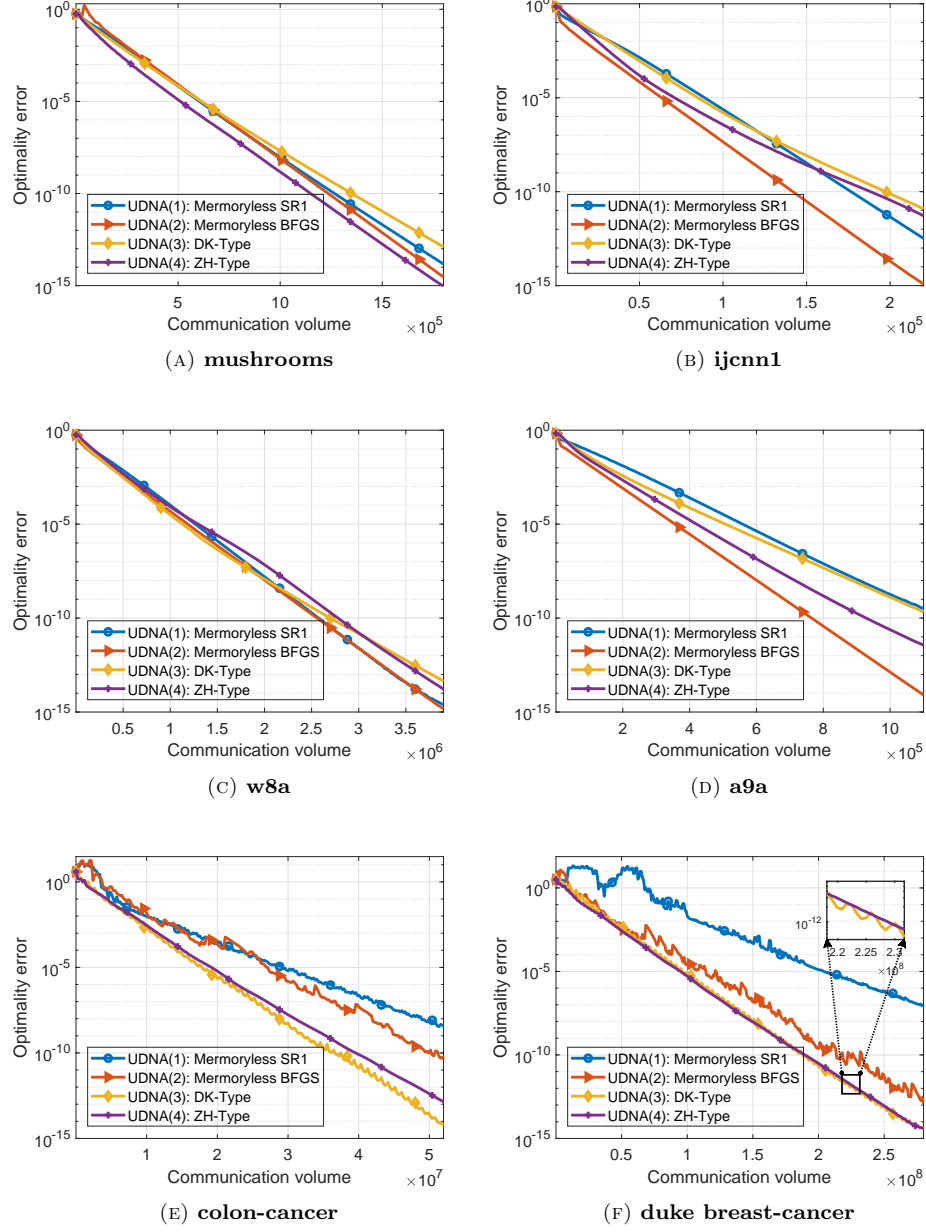
FIGURE 1. Optimality error of UDNAs for minimizing the non-convex logistic regression problem (3.1) on different datasets w.r.t. communication volume.

All algorithms except GUT need two rounds of communication per iteration. GUT needs only one round communication per iteration but uses a decreasing step-size, which yields slow convergence as shown in Fig. 2. Hence, only optimality error
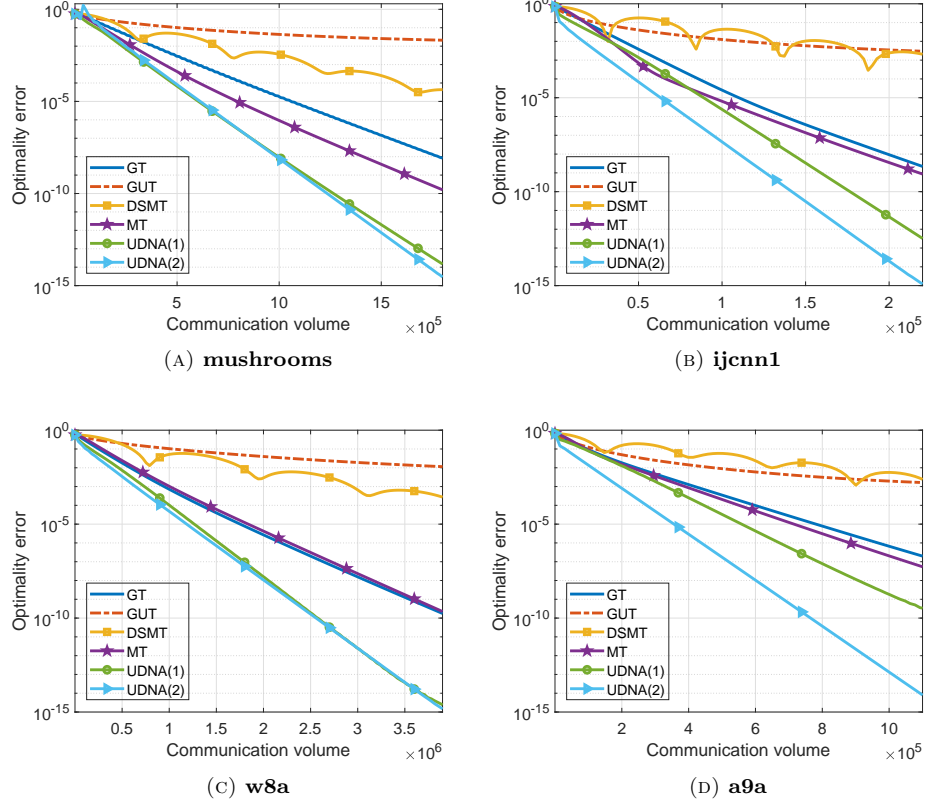
FIGURE 2. Comparisons with gradient-based algorithms for minimizing the nonconvex logistic regression problem (3.1) on different datasets w.r.t. communication volume.

against communication volume is shown in Fig. 2, where we can see our UDNA(1) and UDNA(2) are significantly better than GT and momentum-based methods, including MT and DSMT, for this set of nonconvex classification problems.

We now compare UDNAs with the three well-developed decentralized quasi-Newton algorithms: DQN [47], DR-LM-DFP [62], and D-LM-BFGS [62]. For datasets **mushrooms** (**ijcnn1**; **w8a**; **a9a**), we set $\alpha_i^t = a_* \zeta_i^t$ and $\gamma = 3(3; 2; 3)$ with $a_* = 0.28(0.29; 0.29; 0.3)$ for DQN, where $\zeta_i^t \sim U(0.5, 1.5)$ is the random variable generated over node $i$ at iteration $t$ and satisfies the uniform distribution over interval $(0.5, 1.5)$. For D-LM-BFGS, we set $\alpha = 0.42(0.34; 0.28; 0.36)$, $\epsilon = 10^{-3}$, $\beta = 10^{-3}$, $\mathcal{B} = 10^4$, $\tilde{L} = 5$, $M = 4(4; 5; 8)$. For DR-LM-DFP, we set $\alpha = 0.11(0.18; 0.16; 0.16)$, $\rho = 0.05(0.04; 0.08; 0.05)$, $\epsilon = 10^{-3}$, $\beta = 1(10^{-3}; 10^{-3}; 10^{-3})$, $\mathcal{B} = 10^4$, $\tilde{L} = 5$, $M = 5(5; 4; 5)$.

We see from Fig. 3 that DQN performs best in terms of iteration number for **mushrooms** and **w8a** datasets while UDNA(2) decays fastest for **ijcnn1** and **a9a** datasets. However, as shown in Fig, 4, both UDNA(1) and UDNA(2) are significantly efficient in terms of communication volume when compared with other

algorithms. Despite competitive performance in iteration number, DQN demonstrates significant communication limitations since three rounds of communication are needed in each DQN's iteration.

Finally, we would like to investigate the time efficiency of UDNAs with other advanced algorithms since saving computational overhead is an important advantage of our algorithm. It should be noticed that the computational cost of running algorithms on large sample datasets (**mushrooms**, **ijcnn1**, **w8a**, and **a9a**) is mainly concentrated in computing gradients. Thus, in this part, we focus on **colon-cancer** and **duke breast-cancer** datasets, where the high-dimensional features of the sample pose computational challenges. We select GT, MT, DQN, and D-LM-BFGS as the comparison algorithms. Compared with DR-LM-DFP, one advantage of D-LM-BFGS is that its update can be realized by a two-loop recursion, where the approximate matrix is not generated explicitly and only its multiplications with vectors are computed. Faced with high-dimensional features, one cannot bear the computational cost of DR-LM-DFP. For **colon-cancer** (**duke breast-cancer**) datasets, we set $\eta = 0.0027(0.0013)$ in GT; set $\eta = 0.0035(0.0017)$ and $\beta = 0.8(0.89)$ in MT; set $\alpha_i^t = a_* \zeta_i^t$ and $\gamma = 2$ with $a_* = 0.12(0.1)$ in DQN, where $\zeta_i^t \sim \mathrm{U}(0.5, 1.5)$; set $\alpha = 0.44(0.32)$, $\epsilon = 10^{-3}$, $\beta = 10^{-3}$, $\mathcal{B} = 10^4$, $\tilde{L} = 8(10)$, $M = 13(10)$ in D-LM-BFGS.

From Figs. 5 and 6, we see our UDNAs are the most time-efficient. DQN decays fastest in terms of iteration number but requires a significant amount of computation since matrix-vector products are computed.

## 4. CONCLUSIONS

In this paper, we investigate the unification and generalization of various important algorithms for decentralized optimization problems where $n$ nodes collaboratively minimize the sum of their local nonconvex objectives. We provide a unified algorithmic framework that subsumes a wide range of gradient tracking methods and quasi-Newton methods. Within a cohesive analytical structure, we first show subsequence convergence under mild assumptions with the proper parameter choices and then establish the whole sequence convergence along with a non-asymptotic rate by using the KŁ property of the aggregated objective function. We proceed to generalize this framework by introducing several noval quasi-Newton variants within the unified algorithmic framework. These methods enjoy a very economical computational cost. Our experimental results on the nonconvex decentralized binary classification problem indicate that our newly proposed quasi-Newton algorithms generally exhibit superior performance versus other comparison algorithms, e.g. GT, GUT, DSMT, MT, DQN, D-LM-BFGS, and DR-LM-DFP. Possible future work directions include extensions to stochastic optimization and nonsmooth optimization.

## References

1. Sai Aparna Aketi, Abolfazl Hashemi, and Kaushik Roy, *Global update tracking: A decentralized learning algorithm for heterogeneous data*, Advances in Neural Information Processing Systems, vol. 36, 2023, pp. 48939–48961.
2. Sulaiman A Alghunaim, *Local exact-diffusion for decentralized optimization and learning*, IEEE Transactions on Automatic Control **69** (2024), no. 11, 7371–7386.
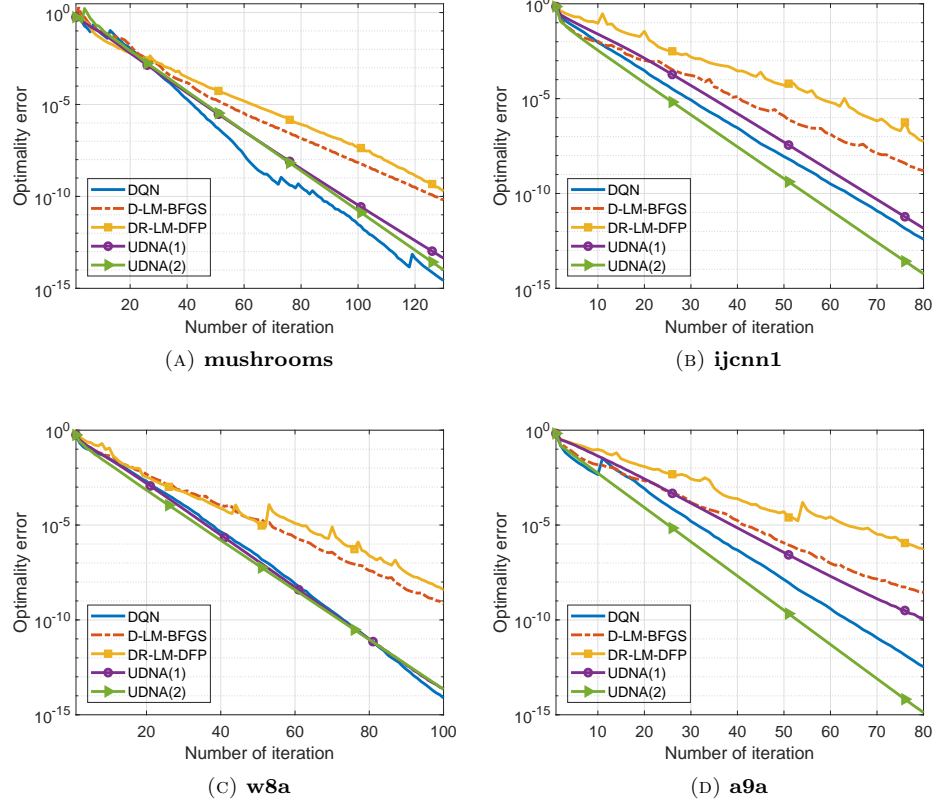
FIGURE 3. Comparisons with quasi-Newton algorithms for minimizing the nonconvex logistic regression problem (3.1) on different datasets w.r.t. number of iteration.

3. Sulaiman A Alghunaim, Ernest K Ryu, Kun Yuan, and Ali H Sayed, *Decentralized proximal gradient algorithms with linear convergence rates*, IEEE Transactions on Automatic Control **66** (2020), no. 6, 2787–2794.

4. Sulaiman A Alghunaim and Kun Yuan, *A unified and refined convergence analysis for nonconvex decentralized learning*, IEEE Transactions on Signal Processing **70** (2022), 3264–3279.

5. Neculai Andrei, *A note on memory-less sr1 and memory-less bfgs methods for large-scale unconstrained optimization*, Numerical Algorithms **90** (2022), no. 1, 223–240.

6. Hedy Attouch and Jérôme Bolte, *On the convergence of the proximal algorithm for nonsmooth functions involving analytic features*, Mathematical Programming **116** (2009), 5–16.

7. Hédy Attouch, Jérôme Bolte, Patrick Redont, and Antoine Soubeyran, *Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the kurdyka-łojasiewicz inequality*, Mathematics of operations research **35** (2010), no. 2, 438–457.

8. Albert S Berahas, Raghu Bollapragada, and Shagun Gupta, *Balancing communication and computation in gradient tracking algorithms for decentralized optimization*, Journal of Optimization Theory and Applications (2024), 1–34.

9. Huiming Chen, Ho-Chun Wu, Shing-Chow Chan, and Wong-Hing Lam, *A stochastic quasi-newton method for large-scale nonconvex optimization with applications*, IEEE transactions on Neural Networks and Learning Systems **31** (2019), no. 11, 4776–4790.

10. Xiaokai Chen, Tianyu Cao, and Gesualdo Scutari, *Enhancing convergence of decentralized gradient tracking under the kl property*, arXiv preprint arXiv:2412.09556 (2024).
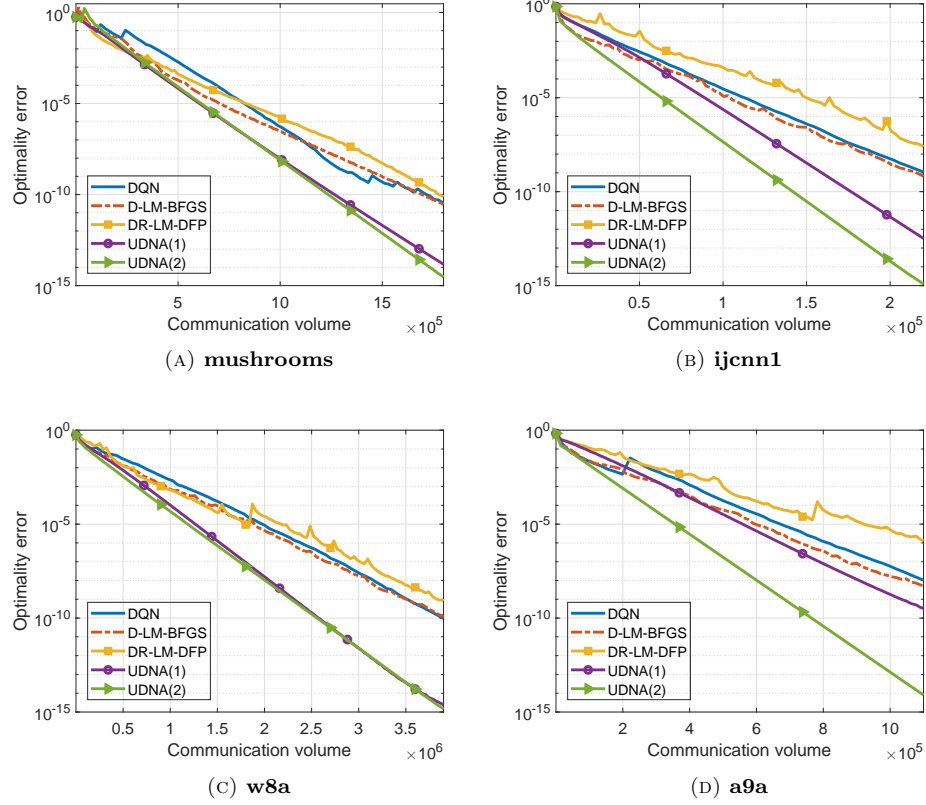
FIGURE 4. Comparisons with quasi-Newton algorithms for minimizing the nonconvex logistic regression problem (3.1) on different datasets w.r.t. communication volume.

11. Frank Curtis, *A self-correcting variable-metric algorithm for stochastic optimization*, International Conference on Machine Learning, PMLR, 2016, pp. 632–641.

12. Yuhong Dai and Caixia Kou, *A nonlinear conjugate gradient algorithm with an optimal property and an improved wolfe line search*, SIAM Journal on Optimization **23** (2013), no. 1, 296–320.

13. Amir Daneshmand, Gesualdo Scutari, and Vyacheslav Kungurtsev, *Second-order guarantees of distributed gradient algorithms*, SIAM Journal on Optimization **30** (2020), no. 4, 3029–3068.

14. Paolo Di Lorenzo and Gesualdo Scutari, *Next: In-network nonconvex optimization*, IEEE Transactions on Signal and Information Processing over Networks **2** (2016), no. 2, 120–136.

15. Haizhou Du, Chaoqian Cheng, and Chengdong Ni, *A unified momentum-based paradigm of decentralized sgd for non-convex models and heterogeneous data*, Artificial Intelligence **332** (2024), 104130.

16. Mark Eisen, Aryan Mokhtari, and Alejandro Ribeiro, *Decentralized quasi-newton methods*, IEEE Transactions on Signal Processing **65** (2017), no. 10, 2613–2628.

17. Giuseppe Fusco and Mario Russo, *A decentralized approach for voltage control by multiple distributed energy resources*, IEEE Transactions on Smart Grid **12** (2021), no. 4, 3115–3127.

18. Juan Gao, Xinwei Liu, Yuhong Dai, Yakui Huang, and Peng Yang, *Achieving geometric convergence for distributed optimization with barzilai-borwein step sizes*, Science China. Information Sciences **65** (2022), no. 4, 149204.
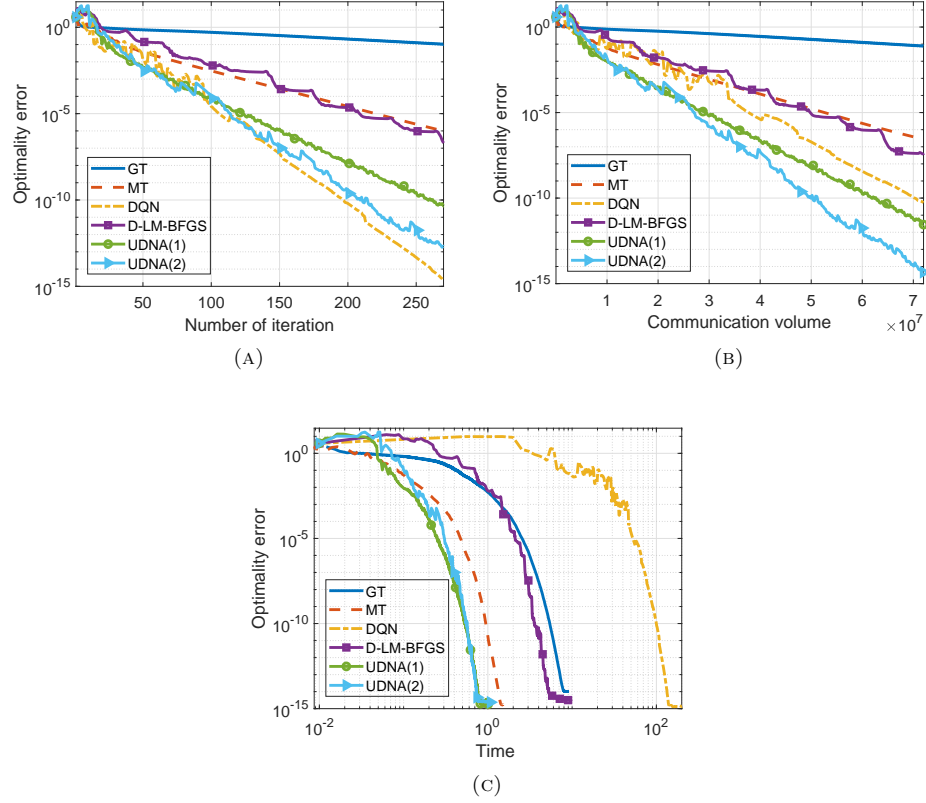
FIGURE 5. Comparisons with quasi-Newton algorithms for minimizing the nonconvex logistic regression problem (3.1) on **colon-cancer** datasets.

19. William W Hager and Hongchao Zhang, *A new conjugate gradient method with guaranteed descent and an efficient line search*, SIAM Journal on optimization **16** (2005), no. 1, 170–192.

20. Kun Huang, Shi Pu, and Angelia Nedić, *An accelerated distributed stochastic gradient method with momentum*, Mathematical Programming (2025), 1–44.

21. Umair Hussan, Huaizhi Wang, Muhammad Ahsan Ayub, Hamna Rasheed, Muhammad Asghar Majeed, Jianchun Peng, and Hui Jiang, *Decentralized stochastic recursive gradient method for fully decentralized opf in multi-area power systems*, Mathematics **12** (2024), no. 19, 3064.

22. Dušan Jakovetić, *A unification and generalization of exact distributed first-order methods*, IEEE Transactions on Signal and Information Processing over Networks **5** (2018), no. 1, 31–46.

23. Dušan Jakovetić, Dragana Bajović, João Xavier, and José MF Moura, *Primal–dual methods for large-scale and distributed convex optimization and data analytics*, Proceedings of the IEEE **108** (2020), no. 11, 1923–1938.

24. Dušan Jakovetić, Nataša Krejić, and Nataša Krklec Jerinkić, *A hessian inversion-free exact second order method for distributed consensus optimization*, IEEE Transactions on Signal and Information Processing over Networks **8** (2022), 755–770.

25. Dušan Jakovetić, Nataša Krejić, and Nataša Krklec Jerinkić, *Exact spectral-like gradient method for distributed optimization*, Computational Optimization and Applications **74** (2019), no. 3, 703–728.
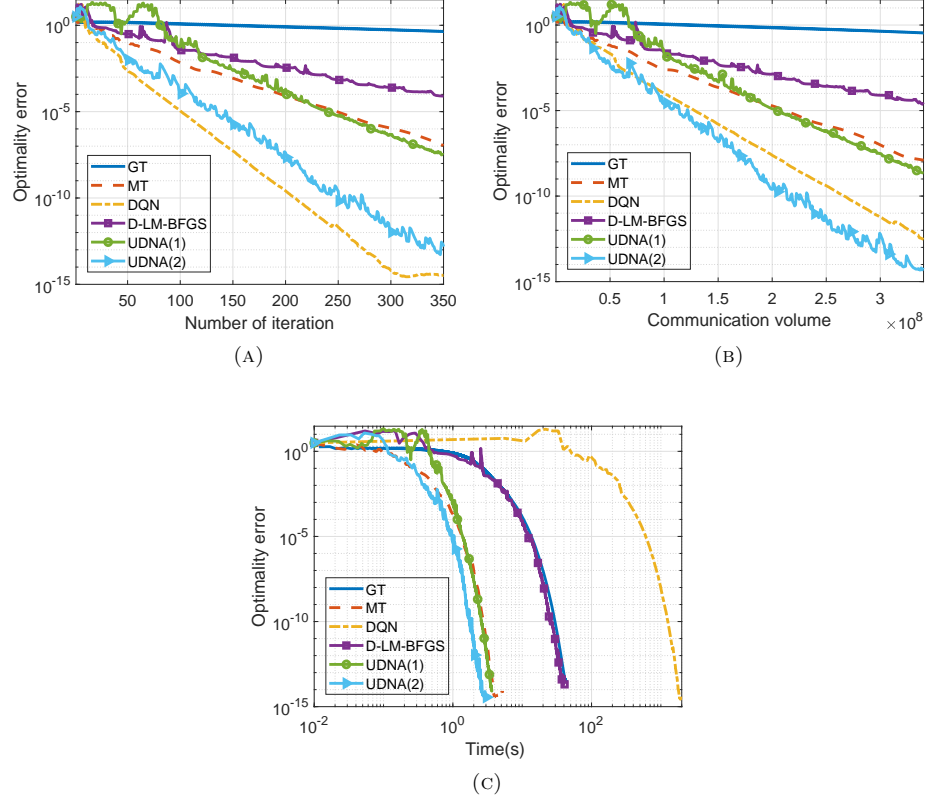
Figure 6. Comparisons with quasi-Newton algorithms for minimizing the nonconvex logistic regression problem (3.1) on **duke breast-cancer** datasets.

26. Eunjeong Jeong, Matteo Zecchin, and Marios Kountouris, *Asynchronous decentralized learning over unreliable wireless networks*, 2022 International Conference on Communications (ICC), 2022, pp. 607–612.

27. Donghui Li and Masao Fukushima, *A modified bfgs method and its global convergence in nonconvex minimization*, Journal of Computational and Applied Mathematics **129** (2001), no. 1-2, 15–35.

28. Guoyin Li and Ting Kei Pong, *Calculus of the exponent of kurdyka–łojasiewicz inequality and its applications to linear convergence of first-order methods*, Foundations of computational mathematics **18** (2018), no. 5, 1199–1232.

29. Yichuan Li, Yonghai Gong, Nikolaos M Freris, Petros Voulgaris, and Dušan Stipanović, *Bfgs-admm for large-scale distributed optimization*, 2021 60th IEEE Conference on Decision and Control (CDC), IEEE, 2021, pp. 1689–1694.

30. Huikang Liu, Jiaojiao Zhang, Anthony Man-Cho So, and Qing Ling, *A communication-efficient decentralized newton's method with provably faster convergence*, IEEE Transactions on Signal and Information Processing over Networks **9** (2023), 427–441.

31. Tao-Wen Liu, *A regularized limited memory bfgs method for nonconvex unconstrained minimization*, Numerical Algorithms **65** (2014), no. 2, 305–323.

32. Stanislaw Lojasiewicz, *Une propriété topologique des sous-ensembles analytiques réels*, Les équations aux dérivées partielles **117** (1963), no. 87-89.

33. Gabriel Mancino-Ball, Yangyang Xu, and Jie Chen, *A decentralized primal-dual framework for non-convex smooth consensus optimization*, IEEE Transactions on Signal Processing **71** (2023), 525–538.

34. Aryan Mokhtari, Qing Ling, and Alejandro Ribeiro, *Network newton distributed optimization methods*, IEEE Transactions on Signal Processing **65** (2016), no. 1, 146–161.

35. Aryan Mokhtari, Wei Shi, Qing Ling, and Alejandro Ribeiro, *A decentralized second-order method with exact linear convergence rate for consensus optimization*, IEEE Transactions on Signal and Information Processing over Networks **2** (2016), no. 4, 507–522.

36. Angelia Nedić, Alex Olshevsky, and Wei Shi, *Achieving geometric convergence for distributed optimization over time-varying graphs*, SIAM Journal on Optimization **27** (2017), no. 4, 2597–2633.

37. Angelia Nedić, Alex Olshevsky, Wei Shi, and César A Uribe, *Geometrically convergent distributed optimization with uncoordinated step-sizes*, 2017 American Control Conference (ACC), IEEE, 2017, pp. 3950–3955.

38. Angelia Nedić and Asuman Ozdaglar, *Distributed subgradient methods for multi-agent optimization*, IEEE Transactions on Automatic Control **54** (2009), no. 1, 48–61.

39. Yitian Qian, Ting Tao, Shaohua Pan, and Houduo Qi, *Convergence of zh-type nonmonotone descent method for kurdyka–łojasiewicz optimization problems*, SIAM Journal on Optimization **35** (2025), no. 2, 1089–1109.

40. Guannan Qu and Na Li, *Harnessing smoothness to accelerate distributed optimization*, IEEE Transactions on Control of Network Systems **5** (2017), no. 3, 1245–1260.

41. Yuben Qu, Haipeng Dai, Yan Zhuang, Jiafa Chen, Chao Dong, Fan Wu, and Song Guo, *Decentralized federated learning for uav networks: Architecture, challenges, and opportunities*, IEEE Network **35** (2022), no. 6, 156–162.

42. Ali H Sayed, *Diffusion adaptation over networks*, Academic Press Library in Signal Processing, vol. 3, 2014, pp. 323–453.

43. Gesualdo Scutari and Ying Sun, *Distributed nonconvex constrained optimization over time-varying digraphs*, Mathematical Programming **176** (2019), 497–544.

44. Wei Shi, Qing Ling, Gang Wu, and Wotao Yin, *Extra: An exact first-order algorithm for decentralized consensus optimization*, SIAM Journal on Optimization **25** (2015), no. 2, 944–966.

45. Wei Shi, Qing Ling, Kun Yuan, Gang Wu, and Wotao Yin, *On the linear convergence of the admm in decentralized consensus optimization*, IEEE Transactions on Signal Processing **62** (2014), no. 7, 1750–1761.

46. Ola Shorinwa and Mac Schwager, *Distributed quasi-newton method for multi-agent optimization*, IEEE Transactions on Signal Processing **72** (2024), 3535–3546.

47. _____, *Distributed quasi-newton method for multi-agent optimization*, IEEE Transactions on Signal Processing **72** (2024), 3535–3546.

48. Zhuoqing Song, Lei Shi, Shi Pu, and Ming Yan, *Optimal gradient tracking for decentralized optimization*, Mathematical Programming **207** (2024), no. 1, 1–53.

49. Akhil Sundararajan, Bryan Van Scoy, and Laurent Lessard, *Analysis and design of first-order distributed optimization algorithms over time-varying graphs*, IEEE Transactions on Control of Network Systems **7** (2020), no. 4, 1597–1608.

50. Yuki Takezawa, Han Bao, Kenta Niwa, Ryoma Sato, and Makoto Yamada, *Momentum tracking: Momentum acceleration for decentralized deep learning on heterogeneous data*, Transactions on Machine Learning Research (2023).

51. Lei Wang, Nachuan Xiao, and Xin Liu, *A double tracking method for optimization with decentralized generalized orthogonality constraints*, arXiv preprint arXiv:2409.04998 (2024).

52. Liping Wang, Hao Wu, and Hongchao Zhang, *A decentralized primal-dual method with quasi-newton tracking*, IEEE Transactions on Signal Processing **73** (2025), 1323–1336.

53. Mou Wu, Haibin Liao, Zhengtao Ding, and Yonggang Xiao, *Music: Accelerated convergence for distributed optimization with inexact and exact methods*, IEEE Transactions on Neural Networks and Learning Systems **36** (2025), no. 3, 4893–4907.

54. Lin Xiao, Stephen Boyd, and Seung-Jean Kim, *Distributed average consensus with least-mean-square deviation*, Journal of Parallel and Distributed Computing **67** (2007), no. 1, 33–46.

55. Ran Xin and Usman A Khan, *Distributed heavy-ball: A generalization and acceleration of first-order methods with gradient tracking*, IEEE Transactions on Automatic Control **65** (2019), no. 6, 2627–2633.

56. Jinming Xu, Ye Tian, Ying Sun, and Gesualdo Scutari, *Distributed algorithms for composite optimization: Unified framework and convergence analysis*, IEEE Transactions on Signal Processing **69** (2021), 3555–3570.

57. Jinming Xu, Shanying Zhu, Yeng Chai Soh, and Lihua Xie, *Augmented distributed gradient methods for multi-agent optimization under uncoordinated constant stepsizes*, 2015 IEEE Conference on Decision and Control (CDC), 2015, pp. 2055–2060.

58. Kun Yuan, Qing Ling, and Wotao Yin, *On the convergence of decentralized gradient descent*, SIAM Journal on Optimization **26** (2016), no. 3, 1835–1854.

59. Kun Yuan, Bicheng Ying, Xiaochuan Zhao, and Ali H Sayed, *Exact diffusion for distributed optimization and learning—part i: Algorithm development*, IEEE Transactions on Signal Processing **67** (2018), no. 3, 708–723.

60. Jinshan Zeng and Wotao Yin, *On nonconvex decentralized gradient descent*, IEEE Transactions on Signal Processing **66** (2018), no. 11, 2834–2848.

61. Jiaojiao Zhang, Qing Ling, and Anthony Man-Cho So, *A newton tracking algorithm with exact linear convergence for decentralized consensus optimization*, IEEE Transactions on Signal and Information Processing over Networks **7** (2021), 346–358.

62. Jiaojiao Zhang, Huikang Liu, Anthony Man-Cho So, and Qing Ling, *Variance-reduced stochastic quasi-newton methods for decentralized learning*, IEEE Transactions on Signal Processing **71** (2023), 311–326.

63. Jiaqi Zhang, Keyou You, and Kai Cai, *Distributed dual gradient tracking for resource allocation in unbalanced networks*, IEEE Transactions on Signal Processing **68** (2020), 2186–2198.

64. Xianyang Zhang, Chen Hu, Bing He, and Zhiguo Han, *Distributed reptile algorithm for meta-learning over multi-agent systems*, IEEE Transactions on Signal Processing **70** (2022), 5443–5456.

## Appendix A. Analytical tools

**Lemma A.1.** *(Young's inequality) For any two vectors* $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^p$,

$$2\mathbf{v}_1^\mathsf{T} \mathbf{v}_2 \leq \eta \|\mathbf{v}_1\|^2 + \frac{1}{\eta}\|\mathbf{v}_2\|^2,$$

$$\|\mathbf{v}_1 + \mathbf{v}_2\|^2 \leq (1+\eta)\|\mathbf{v}_1\|^2 + \left(1 + \frac{1}{\eta}\right)\|\mathbf{v}_2\|^2.$$

**Lemma A.2.** *(Jensen's inequality) For any set of vectors* $\{\mathbf{v}_i\}_{i=1}^n \subset \mathbb{R}^p$,

$$\left\|\frac{1}{n}\sum_{i=1}^n \mathbf{v}_i\right\|^2 \leq \frac{1}{n}\sum_{i=1}^n \|\mathbf{v}_i\|^2.$$

**Lemma A.3.** [10, Lemma 5] *For* $0 \leq a_2 \leq a_1$ *and* $\theta < 1$, *we have*

$$a_1 - a_2 \leq \frac{1}{1-\theta} a_1^\theta \left(a_1^{1-\theta} - a_2^{1-\theta}\right).$$

**Lemma A.4.** [6, Theorem 2] *For the nonnegative sequence* $\{a^t\}$ *satisfying*

$$(a^t)^{\frac{\theta}{1-\theta}} \leq c(a^{t-1} - a^t),$$

*where* $\theta \in (1/2, 1)$ *and* $c > 0$, *there exists* $c' > 0$ *such that*

$$a^t \leq c' t^{-\frac{1-\theta}{2\theta-1}}.$$

---

**Algorithm 2** The memoryless SR1 method –UDNA(1)

---

**Input:** Initial point $\mathbf{x}_i^0$, Maximum iteration T, Stepsize $\alpha > 0$, Mixing matrix $\tilde{\mathbf{W}}$, Parameters $\bar{u} \gg \bar{l} > 0$.

1: Set $t = 0$ and $\mathbf{d}_i^0 = -\mathbf{v}_i^0 = -\mathbf{g}_i^0$.
2: If $t \geq T$, stop.
3: $\mathbf{x}_i^{t+1} = \sum_{j \in \mathcal{N}_i} \tilde{W}_{ij}(\mathbf{x}_j^t + \alpha \mathbf{d}_j^t)$.
4: $\mathbf{v}_i^{t+1} = \sum_{j \in \mathcal{N}_i} \tilde{W}_{ij}(\mathbf{v}_j^t + \mathbf{g}_j^{t+1} - \mathbf{g}_j^t)$.
5:

$$\mathbf{d}_i^{t+1} = \begin{cases} -\mathbf{v}_i^{t+1} - q_i^t(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t), & \text{if } [\lambda_{\min}(\hat{\mathbf{H}}_i^{t+1}), \lambda_{\max}(\hat{\mathbf{H}}_i^{t+1})] \subset [\bar{l}, \bar{u}] \\ & \text{and } (\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)^\mathsf{T} \check{\mathbf{y}}_i^t \neq 0; \\ -\mathbf{v}_i^{t+1}, & \text{otherwise,} \end{cases}$$

where $q_i^t = \frac{(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)^\mathsf{T} \mathbf{v}_i^{t+1}}{(\mathbf{s}_i^t - \check{\mathbf{y}}_i^t)^\mathsf{T} \check{\mathbf{y}}_i^t}$, $\mathbf{s}_i^t = \mathbf{x}_i^{t+1} - \mathbf{x}_i^t$, and $\check{\mathbf{y}}_i^t = \mathbf{v}_i^{t+1} - \mathbf{v}_i^t$.

6: Set $t = t + 1$ and go to Step 2.

**Output:** $\mathbf{x}^T$.

---

**Algorithm 3** The memoryless BFGS method–UDNA(2)

---

**Input:** Initial point $\mathbf{x}_i^0$, Maximum iteration T, Stepsize $\alpha > 0$, Mixing matrix $\tilde{\mathbf{W}}$, Parameters $u \gg l > 0$, $\varrho > 0$.

1: Set $t = 0$ and $\mathbf{d}_i^0 = -\mathbf{v}_i^0 = -\mathbf{g}_i^0$.
2: If $t \geq T$, stop.
3: $\mathbf{x}_i^{t+1} = \sum_{j \in \mathcal{N}_i} \tilde{W}_{ij}(\mathbf{x}_j^t + \alpha \mathbf{d}_j^t)$.
4: $\mathbf{v}_i^{t+1} = \sum_{j \in \mathcal{N}_i} \tilde{W}_{ij}(\mathbf{v}_j^t + \mathbf{g}_j^{t+1} - \mathbf{g}_j^t)$.
5:

$$\mathbf{y}_i^t = \begin{cases} \check{\mathbf{y}}_i^t, & \text{if } (\mathbf{s}_i^t)^\mathsf{T} \check{\mathbf{y}}_i^t > 0 \text{ and } [\lambda_{\min}(H_i^t(\check{\mathbf{y}}_i^t)), \lambda_{\max}(H_i^t(\check{\mathbf{y}}_i^t))] \subset [l, u]; \\ \hat{\mathbf{y}}_i^t, & \text{otherwise,} \end{cases}$$

where

$$\mathbf{s}_i^t = \mathbf{x}_i^{t+1} - \mathbf{x}_i^t, \ \check{\mathbf{y}}_i^t = \mathbf{v}_i^{t+1} - \mathbf{v}_i^t,$$

$$\hat{\mathbf{y}}_i^t = \mathbf{g}_i^{t+1} - \mathbf{g}_i^t + h_i^t \mathbf{s}_i^t, \ h_i^t = \varrho + \max\left\{-\frac{(\mathbf{s}_i^t)^\mathsf{T}(\mathbf{g}_i^{t+1} - \mathbf{g}_i^t)}{\|\mathbf{s}_i^t\|^2}, 0\right\}.$$

6: $\mathbf{d}_i^{t+1} = -\frac{(\mathbf{s}_i^t)^\mathsf{T} \mathbf{y}_i^t}{\|\mathbf{y}_i^t\|^2} \mathbf{v}_i^{t+1} + \left(\frac{(\mathbf{y}_i^t)^\mathsf{T} \mathbf{v}_i^{t+1}}{\|\mathbf{y}_i^t\|^2} - 2\frac{(\mathbf{s}_i^t)^\mathsf{T} \mathbf{v}_i^{t+1}}{(\mathbf{s}_i^t)^\mathsf{T} \mathbf{y}_i^t}\right) \mathbf{s}_i^t + \frac{(\mathbf{s}_i^t)^\mathsf{T} \mathbf{v}_i^{t+1}}{\|\mathbf{y}_i^t\|^2} \mathbf{y}_i^t$.
7: Set $t = t + 1$ and go to Step 2.

**Output:** $\mathbf{x}^T$.

---

APPENDIX B. ALGORITHMIC PROCEDURE FOR UDNAS

School of Mathematics, Nanjing University of Aeronautics and Astronautics
*Email address*: `wuhoo104@nuaa.edu.cn`

School of Mathematics, Nanjing University of Aeronautics and Astronautics
*Email address*: `wlpmath@nuaa.edu.cn`

Department of Mathematics, Louisiana State University
*Email address*: `hozhang@math.lsu.edu`

---

**Algorithm 4** The corrected quasi-Newton method –UDNA(3), UDNA(4)

---

**Input:** Initial point $\mathbf{x}_i^0$, Maximum iteration T, Stepsize $\alpha > 0$, Mixing matrix $\tilde{\mathbf{W}}$,
    Parameters $1 > \lambda > 0$, $\hat{L} > 0$, $\tau \in \{1, 2\}$.
 1: Set $t = 0$ and $\mathbf{d}_i^0 = -\mathbf{v}_i^0 = -\mathbf{g}_i^0$.
 2: If $t \geq T$, stop.
 3: $\mathbf{x}_i^{t+1} = \sum_{j \in \mathcal{N}_i} \tilde{W}_{ij}(\mathbf{x}_j^t + \alpha \mathbf{d}_j^t)$.
 4: $\mathbf{v}_i^{t+1} = \sum_{j \in \mathcal{N}_i} \tilde{W}_{ij}(\mathbf{v}_j^t + \mathbf{g}_j^{t+1} - \mathbf{g}_j^t)$.
 5: $\breve{\mathbf{y}}_i^t = \eta_i^t \check{\mathbf{y}}_i^t + (1 - \eta_i^t)\mathbf{s}_i^t$, where

$$\mathbf{s}_i^t = \mathbf{x}_i^{t+1} - \mathbf{x}_i^t, \; \check{\mathbf{y}}_i^t = \mathbf{v}_i^{t+1} - \mathbf{v}_i^t, \; \eta_i^t = \min\left\{ \hat{\eta}_i^t, \frac{\hat{L}\|\mathbf{s}_i^t\|}{\|\check{\mathbf{y}}_i^t\|} \right\},$$

$$\hat{\eta}_i^t = \left\{ \begin{array}{ll} \frac{(1-\lambda)\|\mathbf{s}_i^t\|^2}{\|\mathbf{s}_i^t\|^2 - (\mathbf{s}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t}, & \text{if } (\mathbf{s}_i^t)^\mathsf{T}\check{\mathbf{y}}_i^t \leq \lambda\|\mathbf{s}_i^t\|^2; \\ 1, & \text{otherwise,} \end{array} \right.$$

 6: $\mathbf{d}_i^{t+1} = -\mathbf{v}_i^{t+1} + \frac{(\mathbf{z}_i^t)^\mathsf{T}\mathbf{v}_i^{t+1}}{2(\mathbf{s}_i^t)^\mathsf{T}\breve{\mathbf{y}}_i^t}\mathbf{s}_i^t + \frac{(\mathbf{s}_i^t)^\mathsf{T}\mathbf{v}_i^{t+1}}{2(\mathbf{s}_i^t)^\mathsf{T}\breve{\mathbf{y}}_i^t}\mathbf{z}_i^t$, where

$$\mathbf{z}_i^t = \breve{\mathbf{y}}_i^t - \tau p_i^t \mathbf{s}_i^t, \; p_i^t = \frac{\|\breve{\mathbf{y}}_i^t\|^2}{(\mathbf{s}_i^t)^\mathsf{T}\breve{\mathbf{y}}_i^t}.$$

 7: Set $t = t + 1$ and go to Step 2.
**Output:** $\mathbf{x}^T$.

---