

D²: Decentralized Training over Decentralized Data

Hanlin Tang^{*1}, Xiangru Lian^{†1}, Ming Yan^{‡4}, Ce Zhang^{§2}, and Ji Liu^{¶3,1}

¹Department of Computer Science, University of Rochester

²Department of Computer Science, ETH Zurich

³Tencent AI Lab

⁴Department of Computational Mathematics, Science and Engineering, Michigan State University

April 23, 2018

Abstract

While training a machine learning model using multiple workers, each of which collects data from their own data sources, it would be most useful when the data collected from different workers can be unique and different. Ironically, recent analysis of decentralized parallel stochastic gradient descent (D-PSGD) relies on the assumption that the data hosted on different workers are not too different. In this paper, we ask the question: Can we design a decentralized parallel stochastic gradient descent algorithm that is less sensitive to the data variance across workers?

In this paper, we present D², a novel decentralized parallel stochastic gradient descent algorithm designed for large data variance among workers (imprecisely, “decentralized” data). The core of D² is a variance reduction extension of the standard D-PSGD algorithm, which improves the convergence rate from $O\left(\frac{\sigma}{\sqrt{nT}} + \frac{(n\zeta^2)^{\frac{1}{3}}}{T^{2/3}}\right)$ to $O\left(\frac{\sigma}{\sqrt{nT}}\right)$ where ζ^2 denotes the variance among data on different workers. As a result, D² is robust to data variance among workers. We empirically evaluated D² on image classification tasks where each worker has access to only the data of a limited set of labels, and find that D² significantly outperforms D-PSGD.

1 Introduction

Training machine learning models in a decentralized way has attracted intensive interests recently Colin et al. [2016], Lian et al. [2017a], Yuan et al. [2016]. In the decentralized setting, there is a set of workers, each of which collects data from different data sources. Instead of sending all of their data to a centralized

^{*}htang14@ur.rochester.edu

[†]xiangru@yandex.com

[‡]yanm@math.msu.edu

[§]ce.zhang@inf.ethz.ch

[¶]ji.liu.uwisc@gmail.com

place, these workers only communicate with their *neighbors*. The goal is to get a model that is the same as if all data are collected in a centralized place. Decentralized learning algorithm is important in scenarios in which centralized communication is expensive or not possible, or the underlying communication network has high latency.

For decentralized learning to provide benefit, each user should provides data that is somehow *unique*, i.e., the variance of data collected from different workers are large. However, many recent theoretical results Lian et al. [2017a,b], Nedic and Ozdaglar [2009], Yuan et al. [2016] all assume a bounded data variance across workers — when data hosted on different workers are very different, these approach could converge slowly, both empirically and theoretically. In this paper, we aim at bringing this discrepancy between the current theoretical understanding and the requirements from *some* practical scenarios.

In this paper, we present D^2 , a novel decentralized learning algorithm designed to be robust under high data variance. The structure and technique of D^2 is built upon standard decentralized parallel stochastic gradient descent (D-PSGD), but benefits from an additional variance reduction component. In the D^2 algorithm, each worker stores the stochastic gradient and its local model in last iterate and linearly combines them with the current stochastic gradient and local model. It results in an improved convergence rate over D-PSGD by eliminating the data variation among workers. In particular, the convergence rate is improved from $O\left(\frac{\sigma}{\sqrt{nT}} + \frac{(n\zeta^2)^{\frac{1}{3}}}{T^{2/3}}\right)$ to $O\left(\frac{\sigma}{\sqrt{nT}}\right)$ where ζ^2 is the data variation among all workers, σ^2 is the data variance within each worker, n is the number of workers, and T is the number of iterations. We empirically show D^2 can significantly outperform D-PSGD by training an image classification model where each worker has access to only the data of a limited set of labels.

Throughout this paper, we consider the following decentralized optimization:

$$\min_{\mathbf{x} \in \mathbb{R}^N} f(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n \overbrace{\mathbb{E}_{\xi \sim \mathcal{D}_i} F_i(\mathbf{x}; \xi)}^{=: f_i(\mathbf{x})}, \quad (1)$$

where n is the number of workers and \mathcal{D}_i is the local data distribution for worker i . All workers are connected to form a connected graph. Each worker can only exchange information with its neighbors.

Definitions and notations Throughout this paper, we use following notations and definitions:

- $\|\cdot\|_F$ denotes the Frobenius norm of matrices.
- $\|\cdot\|$ denotes the ℓ_2 norm for vectors and the spectral norm for matrices.
- $\nabla f(\cdot)$ denotes the gradient of a function f .
- f^* denotes the optimal solution of (1).
- $\lambda_i(\cdot)$ denotes the i th largest eigenvalue of a matrix.
- $\mathbf{x}^{(i)}$ denotes the local model of worker i .
- $\nabla F_i(\mathbf{x}^{(i)}; \xi^{(i)})$ denotes a local stochastic gradient of worker i .
- $\mathbf{1} = [1, 1, \dots, 1]^\top \in \mathbb{R}^n$ denotes the all-one vector.

- In order to organize the algorithm more clearly, here we define the concatenation of all local variables, stochastic gradients, and their average respectively:

$$\begin{aligned}
X &:= [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}] \in \mathbb{R}^{N \times n}, \\
\bar{X} &:= X \frac{\mathbf{1}}{n} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}, \\
G(X; \zeta) &:= [\nabla F_1(\mathbf{x}^{(1)}; \zeta^{(1)}), \dots, \nabla F_n(\mathbf{x}^{(n)}; \zeta^{(n)})] \in \mathbb{R}^{N \times n}, \\
\bar{G}(X, \zeta) &:= G(X, \zeta) \frac{\mathbf{1}}{n} = \frac{1}{n} \sum_{i=1}^n \nabla F_i(\mathbf{x}^{(i)}; \zeta^{(i)}), \\
\nabla f(\bar{X}) &:= \sum_{i=1}^n \frac{1}{n} \nabla f_i(\bar{X}), \\
\bar{\nabla} f(X) &:= \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^{(i)}),
\end{aligned}$$

where ζ is the collection of randomly sampled data from all workers

Organization This paper is organized as follows: Section 2 reviews related work about the proposed approach; Section 3 introduces the state-of-the-art decentralized stochastic gradient descent method and its convergence rate; Section 4 introduces the proposed algorithm and its intuition why it can improve the state-of-the-art approach; and Section 5; Section 6 validates the proposed approaches via empirical study; and Section 7 concludes this paper.

2 Related work

In this section, we review the stochastic gradient descent algorithm and its decentralized variants, decentralized algorithms, and previous variance reduction technologies in this section.

Stochastic gradient descent (SGD) The SGD approaches [Ghadimi and Lan, 2013, Moulines and Bach, 2011, Nemirovski et al., 2009] is quite powerful for solving large-scale machine learning problems. It achieves a convergence rate of $O\left(\frac{1}{\sqrt{T}}\right)$. As an implementation of SGD, the *Centralized Parallel Stochastic Gradient Descent (C-PSGD)*, has been widely used in parallel computation. In C-PSGD, a central worker, whose job is to perform the variable updates, is connected to many leaf workers that are used to compute stochastic gradients in parallel. C-PSGD has been applied to many deep learning frameworks, such as CNTK [Seide and Agarwal, 2016], MXNet [Chen et al., 2015], and TensorFlow [Abadi et al., 2016]. The convergence rate of C-PSGD is $O\left(\frac{1}{\sqrt{nT}}\right)$, which shows it can achieve linear speedup with regards to the number of leaf workers.

Decentralized algorithms Centralized algorithms requires a central server to communicate with all other workers [Suresh et al., 2017]. In contrast, decentralized algorithms can work on any connected network and only rely on the information exchange between neighbor workers [Kashyap et al., 2007, Lavaei and Murray, 2012, Nedic et al., 2009].

Decentralized algorithms are especially useful under a network with limited bandwidth or high latency. It is more favorable when data privacy is sensitive. These advantages have led to successful applications. The decentralized approach for multi-task reinforcement learning was studied in Mhamdi et al. [2017], Omidshafiei et al. [2017]. In Colin et al. [2016], a dual based decentralized algorithm was proposed to solve the pairwise function optimization. Shi et al. [2014] and Mokhtari and Ribeiro [2015] analyzed the decentralized version of the ADMM optimization algorithm. An information theoretic approach was used to analyze decentralization in Dobbe et al. [2017]. The decentralized version of (sub-)gradient descent was studied in Nedic and Ozdaglar [2009], Yuan et al. [2016]. Its $O(1/\sqrt{T})$ convergence requires a diminishing stepsize or a constant stepsize that depends on the total number of iterations. This phenomenon happens because of the variance between the data in different workers, which we call “outer variance” to differentiate it from the variance in SGD. Recently, there are several deterministic decentralized optimization algorithms that allows a constant stepsize. For example, EXTRA Shi et al. [2015a] is the first modification of decentralized gradient descent that converges under a constant stepsize. Later this algorithm is extended for problems with the sum of smooth and nonsmooth functions at each node Shi et al. [2015b]. However, the stepsize depends on both the Lipschitz constant of the differentiable function and the network structure. NIDS is the first algorithm that has a constant network independent stepsize Li et al. [2017]. This algorithm was simultaneously proposed by Yuan et al. [2017] for the smooth case only using a different approach. For directed networks, the algorithm DIGing is proposed in Nedić et al. [2017], where two exchanges are needed in each iteration.¹

Decentralized parallel stochastic gradient descent (D-PSGD) The D-PSGD algorithm [Nedic and Ozdaglar, 2009, Ram et al., 2010a,b] requires each worker to compute a stochastic gradient and exchange its local model with neighbors. In Duchi et al. [2012], a dual averaging based method is proposed for solving the constrained decentralized SGD optimization. In Yuan et al. [2016], the convergence rate for D-PSGD was analyzed when the gradient is assumed to be bounded. In Lan et al. [2017], a decentralized primal-dual type method was proposed with a computational complexity of $O(n/\epsilon^2)$ for general convex objectives. Lian et al. [2017a] proved that D-PSGD can admits linear speedup w.r.t. number of workers with a similar convergence rate like C-PSGD.

Variance reduction technology There have been many methods developed for reducing the variance in SGD, including SVRG [Johnson and Zhang, 2013], SAGA [Defazio et al., 2014], SAG [Schmidt et al., 2017], MISO [Mairal, 2015], and mS2GD [Konečný et al., 2016]. However, most of these technologies are just designed for the centralized approaches. The DSA algorithm [Mokhtari and Ribeiro, 2016] applies the variance reduction similar to SAGA on strongly convex decentralized optimization problems and proved a linear convergence rate. However, the speedup property is unclear and a table of all stochastic gradients need to be stoblack.

3 Preliminary: decentralized stochastic gradient descent

The decentralized stochastic gradient descent [Lian et al., 2017a, Shahrampour and Jadbabaie, 2017, Zhang et al., 2017] allows each worker (say worker i) maintaining its own local variable $x^{(i)}$. During each iteration (say, iteration t), each worker performs the following steps:

¹To Prof. Yan: could you write couple of sentences to summarize these papers.

1. Query its neighbors' local variables.
2. Take weighted average with its local variable and neighbors' local variables:

$$\mathbf{x}_{t+\frac{1}{2}}^{(i)} = \sum_{j=1}^n W_{ij} \mathbf{x}_t^{(j)}$$

where W_{ij} is the (i, j) element of the matrix W , $W_{ij} = 0$ means worker i and worker j are not connected.

3. Perform one stochastic gradient descent step

$$\mathbf{x}_{t+1}^{(i)} = \mathbf{x}_{t+\frac{1}{2}}^{(i)} - \gamma \nabla F(\mathbf{x}_t^{(i)}; \xi_t^{(i)})$$

where $\xi_t^{(i)}$ represents the data sampled in worker i at the iteration t following the distribution \mathcal{D}_i .

From a global point of view, the update rule D-PSGD algorithm can be viewed as

$$X_{t+1} = X_t W - \gamma G(X_t; \xi_t).$$

It admits the following rate shown in Theorem 1.

Theorem 1 (Convergence rate of D-PSGD [Lian et al., 2017a]). *Under certain assumptions, the output of D-PSGD admits the following inequality*

$$\frac{1 - \gamma L}{2T} \sum_{t=0}^{T-1} \mathbb{E} \left\| \overline{\nabla f}(X_t) \right\|^2 + \frac{D_1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left\| \nabla f(\bar{X}_t) \right\|^2 \leq \frac{f(0) - f^*}{\gamma T} + \frac{\gamma L}{2n} \sigma^2 + \frac{\gamma^2 L^2 n \sigma^2}{(1 - \lambda) D_2} + \frac{9\gamma^2 L^2 n \zeta^2}{(1 - \sqrt{\lambda})^2 D_2},$$

where ρ reflects the property of the network, D_1 and D_2 are defined to be

$$D_1 := \left(\frac{1}{2} - \frac{9\gamma^2 L^2 n}{(1 - \sqrt{\rho})^2 D_2} \right)$$

$$D_2 := \left(1 - \frac{18\gamma^2}{(1 - \sqrt{\rho})^2} n L^2 \right)$$

and σ and ζ measure the variation within each worker and among all workers respectively

$$\mathbb{E}_{\xi \sim \mathcal{D}_i} \left\| \nabla F_i(\mathbf{x}; \xi) - \nabla f_i(\mathbf{x}) \right\|^2 \leq \sigma^2, \quad \forall i, \forall \mathbf{x}, \quad (2)$$

$$\frac{1}{n} \sum_{i=1}^n \left\| \nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x}) \right\|^2 \leq \zeta^2, \quad \forall i, \forall \mathbf{x}. \quad (3)$$

Choosing the optimal steplength $\gamma = \frac{1}{L + \sigma \sqrt{\frac{K}{n} + n^{\frac{1}{3}} \zeta^{\frac{2}{3}} T^{\frac{1}{3}}}}$ we have the following convergence rate:

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} (\left\| \nabla f(\bar{X}_t) \right\|^2) \leq O \left(\frac{\sigma}{\sqrt{nT}} + \frac{n^{\frac{1}{3}} \zeta^{\frac{2}{3}}}{T^{\frac{2}{3}}} + \frac{1}{T} \right).$$

The proposed D^2 algorithm can improve the convergence rate by removing the dependence to the global bound of outer variance ζ .

Algorithm 1 D^2 algorithm

- 1: **Input:** Initial point $\mathbf{x}_0^{(i)} = \mathbf{0}$, iteration step length γ , confusion matrix W , and the total number of iterations T
- 2: **for** $t = 0, 1, 2, \dots, T$ **do**
- 3: Randomly sample $\zeta_t^{(i)}$ from the local data of the i th worker.
- 4: Compute a local stochastic gradient based on $\zeta_t^{(i)}$ and current optimization variable $\mathbf{x}_t^{(i)}$: $\nabla F_i(\mathbf{x}_t^{(i)}; \zeta_t^{(i)})$.
- 5:
- 6: **if** $t=0$ **then**
- 7: $\mathbf{x}_{t+\frac{1}{2}}^{(i)} = \mathbf{x}_t^{(i)} - \gamma \nabla F_i(\mathbf{x}_t^{(i)}; \zeta_t^{(i)})$,
- 8: **else**
- 9: $\mathbf{x}_{t+\frac{1}{2}}^{(i)} = 2\mathbf{x}_t^{(i)} - \mathbf{x}_{t-1}^{(i)} - \gamma \nabla F_i(\mathbf{x}_t^{(i)}; \zeta_t^{(i)}) + \gamma \nabla F_i(\mathbf{x}_{t-1}^{(i)}; \zeta_{t-1}^{(i)})$.
- 10: **end if**
- 11: Each worker sends $\mathbf{x}_{t+\frac{1}{2}}^{(i)}$ to its neighbors, and take the weighted average

$$\mathbf{x}_{t+1}^{(i)} = \sum_{j=1}^n W_{ij} \mathbf{x}_{t+\frac{1}{2}}^{(j)},$$

where $\mathbf{x}_{t+\frac{1}{2}}^{(j)}$ is from the worker j .

- 12: **end for**
 - 13: **Output:** $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_T^{(i)}$
-

4 The D^2 algorithm

In D^2 algorithm, each worker repeats the following updating rule (say, at iteration t) for worker i

1. Compute a local stochastic gradient $\nabla F(\mathbf{x}_t^{(i)}; \zeta_t^{(i)})$ by sampling $\zeta_t^{(i)}$ from distribution $\mathcal{D}^{(i)}$;
2. Update the local model $\mathbf{x}_{t+\frac{1}{2}}^{(i)} \leftarrow 2\mathbf{x}_t^{(i)} - \mathbf{x}_{t-1}^{(i)} - \gamma \nabla F_i(\mathbf{x}_t^{(i)}; \zeta_t^{(i)}) + \gamma \nabla F_i(\mathbf{x}_{t-1}^{(i)}; \zeta_{t-1}^{(i)})$ using the local models and stochastic gradients in both the t th iteration and the $(t-1)$ th iteration.
3. When the synchronization barrier is met, exchange $\mathbf{x}_{t+\frac{1}{2}}^{(i)}$ with neighbors:

$$\mathbf{x}_{t+1}^{(i)} = \sum_{j=1}^n W_{ij} \mathbf{x}_{t+\frac{1}{2}}^{(j)}.$$

From a global point of view, the update rule of D^2 can be viewed as:

$$X_{t+1} = (2X_t - X_{t-1} - \gamma G(X_t; \zeta_t) + \gamma G(X_{t-1}; \zeta_{t-1})) W.$$

The complete algorithm is summarized in Algorithm 1.

D² essentially runs the stochastic gradient descent step. To understand the intuition of D², let us consider the mean value \bar{X}_t , which gets updated just like the standard stochastic gradient descent:

$$\begin{aligned}\bar{X}_{t+1} &= (2X_t - X_{t-1} - \gamma G(X_t; \xi_t) + \gamma G(X_{t-1}; \xi_{t-1})) W \frac{\mathbf{1}_n}{n}, \\ \bar{X}_{t+1} &= 2\bar{X}_t - \bar{X}_{t-1} - \gamma \bar{G}(X_t; \xi_t) + \gamma \bar{G}(X_{t-1}; \xi_{t-1}),\end{aligned}$$

or equivalently

$$\begin{aligned}\bar{X}_{t+1} - \bar{X}_t &= \bar{X}_t - \bar{X}_{t-1} - \gamma \bar{G}(X_t; \xi_t) + \gamma \bar{G}(X_{t-1}; \xi_{t-1}), \\ &= \bar{X}_1 - \bar{X}_0 - \gamma \sum_{k=1}^t (\bar{G}(X_k; \xi_k) - \bar{G}(X_{k-1}; \xi_{k-1})) \\ &= -\gamma \bar{G}(X_t; \xi_t). \quad (\text{due to } X_1 = X_0 - \gamma G(X_0; \xi_0)).\end{aligned}\tag{4}$$

Why D² improves the D-PSGD? Acute reviewers may notice that the D-PSGD algorithm also essentially updates in the form of stochastic gradient descent in (4). Then why D² can improve D-PSGD?

Assume that X_t has achieved the optimum $X^* := x^* \mathbf{1}^\top$ with all local models equal to the optimum x^* to (1). Then for D-PSGD, the next update will be

$$X_{t+1} = X^* - \gamma G(X^*; \xi_t).$$

It shows that the convergence when we approach a solution is affected by $\mathbb{E}[\|G(X^*; \xi_t)\|_F^2]$, which is bounded by

$$\mathcal{O}(\sigma^2 + \zeta^2).$$

as we can see from the following:

$$\begin{aligned}& \mathbb{E}[\|G(X^*; \xi_t)\|_F^2] \\ &= \mathbb{E} \sum_{i=1}^n \left\| \left(\nabla F_i(\mathbf{x}^*; \xi_{t+1}^{(i)}) - \nabla f_i(\mathbf{x}^*) \right) + \nabla f_i(\mathbf{x}^*) \right\|^2 \\ &\leq 2\mathbb{E} \sum_{i=1}^n \left\| \left(\nabla F_i(\mathbf{x}^*; \xi_{t+1}^{(i)}) - \nabla f_i(\mathbf{x}^*) \right) \right\|^2 + 2 \|\nabla f_i(\mathbf{x}^*) - \nabla f(\mathbf{x}^*)\|^2 \\ &\leq 2\sigma^2 + 2\zeta^2.\end{aligned}$$

Next we apply a similar analysis for D² by assuming that both X_{t-1} and X_t have reached the optimal solution X^* . The next update for D² will be:

$$X_{t+1} = (X^* - \gamma G(X^*; \xi_t) - \gamma G(X^*; \xi_{t-1})) W.$$

It shows that for D², the convergence when we approach a solution relies on the magnitude of $\mathbb{E}[\|G(X^*; \xi_t) - G(X^*; \xi_{t-1})\|_F^2]$, which is bounded by:

$$\mathcal{O}(\sigma^2),$$

which can be seen from:

$$\begin{aligned}\mathbb{E}[\|G(X^*; \xi_t) - G(X^*; \xi_{t-1})\|_F^2] &= \mathbb{E} \sum_{i=1}^n \left\| \nabla F_i(\mathbf{x}^*; \xi_t^{(i)}) - \nabla f_i(\mathbf{x}^*) \right\|^2 - \mathbb{E} \sum_{i=1}^n \left\| \nabla F_i(\mathbf{x}^*; \xi_{t-1}^{(i)}) - \nabla f_i(\mathbf{x}^*) \right\|^2 \\ &\leq 2\sigma^2.\end{aligned}$$

5 Theoretical guarantee

This section provides the theoretical guarantee for the proposed D^2 algorithm. We first give the assumptions requirblack below.

Assumption 1. *Throughout this paper, we make the following commonly used assumptions:*

1. **Lipschitzian gradient:** All function $f_i(\cdot)$'s are with L -Lipschitzian gradients.
2. **Bounded variance:** Assume bounded variance of stochastic gradient within each worker

$$\mathbb{E}_{\xi \sim \mathcal{D}_i} \|\nabla F_i(\mathbf{x}; \xi) - \nabla f_i(\mathbf{x})\|^2 \leq \sigma^2, \quad \forall i, \forall \mathbf{x}.$$

3. **Symmetric confusion matrix:** The confusion matrix W is symmetric and satisfies $W\mathbf{1} = \mathbf{1}$.
4. **Spectral gap:** Let the eigenvalues of $W \in \mathbb{R}^{n \times n}$ be $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. Denote by for short

$$\lambda := \max_{i \in \{2, \dots, n\}} \lambda_i = \lambda_2.$$

We assume $\lambda < 1$ and $\lambda_n > -\frac{1}{3}$.

5. **Initialization:** W.l.o.g., assume all local variables are initialized by zero, that is, $X_0 = 0$.

Existing decentralized consensus algorithms [Li et al., 2017, Shi et al., 2015b] use a modification of the doubly stochastic matrix such that $\lambda > 0$, i.e., choose $W = (\tilde{W} + I)/2$ where W is a doubly stochastic matrix. Recently, Li and Yan [2017] show that $\lambda_n > -1/3$ is optimal in the convergence of EXTRA. However, the optimal λ_n for NIDS [Li et al., 2017] is unknown. In this paper, we proved that $-\frac{1}{3}$ is the infimum of λ_n , and when it blackuces to deterministic case, this condition is weaker than that in Li et al. [2017]. This is important, because we actually can use a W that performs better.

Given Assumption 1, we have following convergence guarantee for D^2 :

Theorem 2 (Convergence of Algorithm 1). *Choose the steplength γ in Algorithm 1 to be a constant satisfying $1 - 24C_2\gamma^2L^2 > 0$. Under Assumption 1, we have the following convergence rate for Algorithm 1:*

$$\begin{aligned} & A_1 \|\nabla f(\mathbf{0})\|^2 + \sum_{t=1}^{T-1} \left(\mathbb{E} \|\nabla f(\bar{X}_t)\|^2 + A_2 \mathbb{E} \|\bar{\nabla} f(X_t)\|^2 \right) \\ & \leq \frac{2(f(0) - f^*)}{\gamma} + \frac{LT\gamma\sigma^2}{n} + \frac{6L^2C_1\gamma^2\zeta_0^2}{C_3} + \frac{12L^2C_2\gamma^2\sigma^2T}{C_3} + \frac{6L^2C_2\gamma^4L^2\sigma^2T}{nC_3} + \frac{6L^2C_1\gamma^2\sigma^2}{C_3}, \end{aligned} \quad (5)$$

where

$$\begin{aligned} \zeta_0 &:= \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{0}) - \nabla f(\mathbf{0})\|^2, \\ v &:= \lambda_n - \sqrt{\lambda_n^2 - \lambda_n}, \\ C_1 &:= \max \left\{ \frac{1}{1 - |v|^2}, \frac{1}{(1 - \lambda)^2} \right\} \geq 1, \end{aligned}$$

$$\begin{aligned}
C_2 &:= \max \left\{ \frac{\lambda_n^2}{(1 - |\nu|^2)}, \frac{\lambda^2}{(1 - \sqrt{\lambda})^2(1 - \lambda)} \right\}, \\
C_3 &:= 1 - 24C_2\gamma^2L^2, \\
A_1 &:= 1 - \frac{6L^2C_1\gamma^2}{C_3}, \\
A_2 &:= 1 - L\gamma - \frac{6L^2C_2\gamma^4L^2}{C_3}.
\end{aligned}$$

By appropriately specifying the step length γ we reach the following corollary:

Corollary 3. Choose the step length γ in Algorithm 1 to be $\gamma = \frac{1}{8\sqrt{C_2}L + 6\sqrt{C_1}L + \sigma\sqrt{\frac{T}{n}}}$, where C_1 and C_2 are defined in Theorem 2. Under Assumption 1, the following convergence rate holds

$$\frac{1}{T} \sum_{t=0}^T \mathbb{E} \|\nabla f(\bar{X}_t)\|^2 \lesssim \frac{\sigma}{\sqrt{nT}} + \frac{1}{T} + \frac{\zeta_0^2}{T + \sigma^2 T^2} + \frac{\sigma^2}{1 + \sigma^2 T},$$

where ζ_0 is defined in Theorem 2 and we treat $f(0) - f^*$, L , λ_n , and λ as constants.

Note that we can obtain even better constants by choosing different parameters and applying tighter inequalities, however, the main result of this corollary is to show the order of the convergence. We highlight a few key observations from our theoretical results in the following.

Tightness of the convergence rate Setting $\sigma = 0$ and $\zeta_0 = 0$, which blackuces the VR-SGD to a normal GD algorithm, we shall see that the convergence rate becomes $O\left(\frac{1}{T}\right)$, which is exactly the rate of GD.

Linear speedup Since the leading term of the convergence rate is $O\left(\frac{1}{\sqrt{nT}}\right)$, which is consistent with the convergence rate of C-PSGD, this indicates that we would achieve a linear speed up with respect to the number of nodes.

Consistent with NIDS In NIDS [Li and Yan, 2017], the term depends on ζ_0 in the convergence rate is $O\left(\frac{\zeta_0^2}{T}\right)$. While the corresponding term in D^2 is $O\left(\frac{\zeta_0^2}{T + \sigma^2 T^2}\right)$, which indicates when our algorithm is consistent with NIDS because in NIDS σ is considerblack to be 0.

Superiority over D-PSGD When compablack to D-PSGD, the convergence rate of D^2 only depends on ζ_0 , and the corresponding decaying rate is $\frac{\zeta_0}{T^2}$. Whereas in D-PSGD [Lian et al., 2017a], we need to assume an upper bound for the global variance between different nodes' dataset, and its influence can be compablack to σ^2 , the inner variance of each node itself. This means we can always achieve a much better convergence rate than D-PSGD.

6 Experiments

We evaluate the effectiveness of D^2 by comparing it with both centralized and decentralized SGD algorithms.

6.1 Experiment Settings

We conduct experiments in two settings.

1. **TRANSFERLEARNING**: We test the case that each worker has access to a local pre-trained neural network as feature extractor, and we want to train a logistic regression model among all these workers. In our experiment, we select the first 16 classes of ImageNet and use InceptionV4 as the feature extractor to extract 2048 features for each image. We conduct data augmentation and generate a blurblack version for each image. In total this dataset contains $16 \times 1300 \times 2$ images.
2. **LENET**: We test the case that all workers collaboratively train a neural network model. We train a LeNet on the CIFAR10 dataset. In total this dataset contains 50,000 images of size 32×32 .

One caveat of training more recent neural networks is that modern architectures often have a batch normalization layer, which inherently assumes that the data distribution is uniform across different batches, which is not the case that we are interested in. In principle, we could also flow the batch information through the network in a decentralized way; however, we leave this as future work.

By default, each worker only has *exclusive* access to a subset of classes. For **TRANSFERLEARNING**, we use 16 workers and each worker has access to one class; for **LENET**, we use 5 workers and each worker has access to two classes. For comparison, we also consider a case when the datasets is first shuffled and then uniformly partitioned among all the workers, we call this the *shuffled case*, and the default one the *unshuffled case*. We use a ring topology for both experiments.

Parameter Tuning. For **TRANSFERLEARNING**, we use constant learning rates and tune it from $\{0.01, 0.025, 0.05, 0.075, 0.1\}$. For **LENET**, we use constant learning rate 0.05 which is tuned from $\{0.5, 0.1, 0.05, 0.01\}$ for centralized algorithms and batch size 128 on each worker.

Metrics. In this paper, we mainly focus on the convergence rate of different algorithms instead of the wall clock speed. This is because the implementation of D^2 is a minor change over the standard D-PSGD algorithm, and thus they has almost the same speed to finish one epoch of training, and both are no slower than the centralized algorithm. When the network has high latency, if a decentralized algorithm (D^2 or D-PSGD) converges with a similar speed as the centralized algorithm, it can be up to one order of magnitude faster Lian et al. [2017a]. However, the convergence rate depending on the “outer variance” is different for both algorithms.

6.2 Unshuffled Case

We are mostly interested in the unshuffled case, in which the data variation across workers is maximized. Figure 1 shows the result. In the unshuffled case, we see that the D-PSGD algorithm convergences slower than the centralized case. This is consistent with the original D-PSGD paper [Lian et al., 2017a]. On the other hand, D^2 converges much faster than D-PSGD, and achieves almost the same loss as the centralized algorithm. For the LeNet case, each worker only has access to data of assigned two labels, which means the data variation is very large. The D-PSGD does not converge with the the given learning rate 0.05.²

²We can tune the learning rate 50x smaller for D-PSGD to converge in this case, but doing so will make D-PSGD stuck at the starting point for quite a long time.

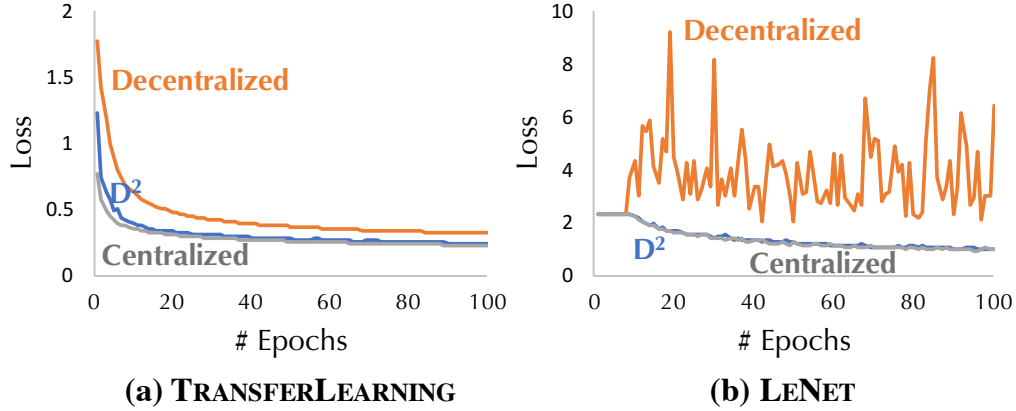


Figure 1: Convergence of Different Distributed Training Algorithms (Unshuffled Case).

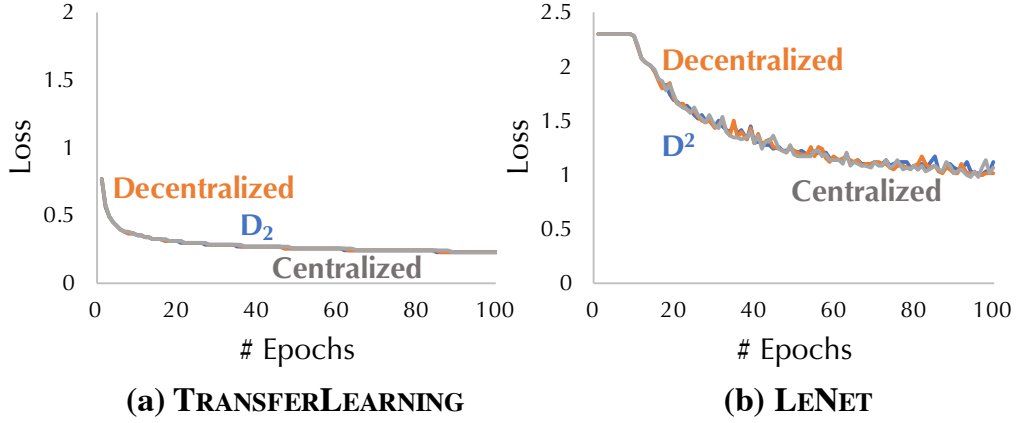


Figure 2: Convergence of Different Distributed Training Algorithms (Shuffled Case).

6.3 Shuffled Case

As a sanity check, Figure 2 shows the result of three different algorithms on the shuffled data. In this case, the data variation of among workers is small (in expectation, they are drawn from the same distribution). We see that, all strategies have similar convergence rate. This validate that the D^2 algorithm is more effective for larger data variation between different workers.

7 Conclusion

In this paper, we propose a decentralized algorithm, namely, D^2 algorithm. D^2 algorithm integrates the D-PSGD algorithm with the variance reduction technology, by which we improves the convergence rate of D-PSGD. The variance reduction technology used in this paper is different from the commonly used ones such as SVRG and SAGA, that are designed for centralized approaches. Experiments validate the advantage of D^2 over D-PSGD — D^2 converges with a rate that is similar to centralized SGD while D-PSGD does not converge to the a solution with a similar quality when the data variance is large. While being robust to large data variance among workers, the same performance benefit of D-PSGD over the centralized strategy still holds for D^2 .

References

- M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. Tensorflow: A system for large-scale machine learning. In *OSDI*, volume 16, pages 265–283, 2016.
- T. Chen, M. Li, Y. Li, M. Lin, N. Wang, M. Wang, T. Xiao, B. Xu, C. Zhang, and Z. Zhang. Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv preprint arXiv:1512.01274*, 2015.
- I. Colin, A. Bellet, J. Salmon, and S. Cl  men  on. Gossip dual averaging for decentralized optimization of pairwise functions. In *International Conference on Machine Learning*, pages 1388–1396, 2016.
- A. Defazio, F. Bach, and S. Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in neural information processing systems*, pages 1646–1654, 2014.
- R. Dobbe, D. Fridovich-Keil, and C. Tomlin. Fully decentralized policies for multi-agent systems: An information theoretic approach. In *Advances in Neural Information Processing Systems*, pages 2945–2954, 2017.
- J. C. Duchi, A. Agarwal, and M. J. Wainwright. Dual averaging for distributed optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic control*, 57(3):592–606, 2012.
- S. Ghadimi and G. Lan. Stochastic first- and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013. doi: 10.1137/120880811.
- R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in neural information processing systems*, pages 315–323, 2013.
- A. Kashyap, T. Ba  sar, and R. Srikant. Quantized consensus. *Automatica*, 43(7):1192–1203, 2007.
- J. Kone  n  , J. Liu, P. Richt  rik, and M. Tak   . Mini-batch semi-stochastic gradient descent in the proximal setting. *IEEE Journal of Selected Topics in Signal Processing*, 10(2):242–255, 2016.
- G. Lan, S. Lee, and Y. Zhou. Communication-efficient algorithms for decentralized and stochastic optimization. 01 2017.
- J. Lavaei and R. M. Murray. Quantized consensus by means of gossip algorithm. *IEEE Transactions on Automatic Control*, 57(1):19–32, 2012.
- Z. Li and M. Yan. A primal-dual algorithm with optimal stepsizes and its application in decentralized consensus optimization. *arXiv preprint arXiv:1711.06785*, 2017.
- Z. Li, W. Shi, and M. Yan. A decentralized proximal-gradient method with network independent step-sizes and separated convergence rates. *arXiv preprint arXiv:1704.07807*, 2017.
- X. Lian, C. Zhang, H. Zhang, C.-J. Hsieh, W. Zhang, and J. Liu. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. 05 2017a.
- X. Lian, W. Zhang, C. Zhang, and J. Liu. Asynchronous decentralized parallel stochastic gradient descent. *arXiv preprint arXiv:1710.06952*, 2017b.
- J. Mairal. Incremental majorization-minimization optimization with application to large-scale machine learning. *SIAM Journal on Optimization*, 25(2):829–855, 2015.

- E. Mhamdi, E. Mahdi, H. Hendrikx, R. Guerraoui, and A. D. O. Maurer. Dynamic safe interruptibility for decentralized multi-agent reinforcement learning. Technical report, EPFL, 2017.
- A. Mokhtari and A. Ribeiro. Decentralized double stochastic averaging gradient. In *Signals, Systems and Computers, 2015 49th Asilomar Conference on*, pages 406–410. IEEE, 2015.
- A. Mokhtari and A. Ribeiro. Dsa: Decentralized double stochastic averaging gradient algorithm. *Journal of Machine Learning Research*, 17(61):1–35, 2016.
- E. Moulines and F. R. Bach. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*, pages 451–459. Curran Associates, Inc., 2011.
- A. Nedic and A. Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.
- A. Nedic, A. Olshevsky, A. Ozdaglar, and J. N. Tsitsiklis. On distributed averaging algorithms and quantization effects. *IEEE Transactions on Automatic Control*, 54(11):2506–2517, 2009.
- A. Nedić, A. Olshevsky, and M. G. Rabbat. Network topology and communication-computation tradeoffs in decentralized optimization. *arXiv preprint arXiv:1709.08765*, 2017.
- A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009. doi: 10.1137/070704277.
- S. Omidshafiei, J. Papis, C. Amato, J. P. How, and J. Vian. Deep decentralized multi-task multi-agent rl under partial observability. *arXiv preprint arXiv:1703.06182*, 2017.
- S. S. Ram, A. Nedić, and V. V. Veeravalli. Asynchronous gossip algorithm for stochastic optimization: Constant stepsize analysis. In *Recent Advances in Optimization and its Applications in Engineering*, pages 51–60. Springer, 2010a.
- S. S. Ram, A. Nedić, and V. V. Veeravalli. Distributed stochastic subgradient projection algorithms for convex optimization. *Journal of optimization theory and applications*, 147(3):516–545, 2010b.
- M. Schmidt, N. Le Roux, and F. Bach. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162(1-2):83–112, 2017.
- F. Seide and A. Agarwal. Cntk: Microsoft’s open-source deep-learning toolkit. In *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’16*, pages 2135–2135, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-4232-2. doi: 10.1145/2939672.2945397.
- S. Shahrampour and A. Jadbabaie. Distributed online optimization in dynamic environments using mirror descent. *IEEE Transactions on Automatic Control*, 2017.
- W. Shi, Q. Ling, K. Yuan, G. Wu, and W. Yin. On the linear convergence of the admm in decentralized consensus optimization. *IEEE Trans. Signal Processing*, 62(7):1750–1761, 2014.
- W. Shi, Q. Ling, G. Wu, and W. Yin. Extra: An exact first-order algorithm for decentralized consensus optimization. *SIAM Journal on Optimization*, 25(2):944–966, 2015a.
- W. Shi, Q. Ling, G. Wu, and W. Yin. A proximal gradient algorithm for decentralized composite optimization. *IEEE Transactions on Signal Processing*, 63(22):6013–6023, 2015b.

- A. T. Suresh, F. X. Yu, S. Kumar, and H. B. McMahan. Distributed mean estimation with limited communication. In D. Precup and Y. W. Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 3329–3337, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR. URL <http://proceedings.mlr.press/v70/suresh17a.html>.
- K. Yuan, Q. Ling, and W. Yin. On the convergence of decentralized gradient descent. *SIAM Journal on Optimization*, 26(3):1835–1854, 2016. doi: 10.1137/130943170.
- K. Yuan, B. Ying, X. Zhao, and A. H. Sayed. Exact diffusion for distributed optimization and learning—part i: Algorithm development. *arXiv preprint arXiv:1702.05122*, 2017.
- W. Zhang, P. Zhao, W. Zhu, S. C. Hoi, and T. Zhang. Projection-free distributed online learning in networks. In *International Conference on Machine Learning*, pages 4054–4062, 2017.

Supplemental Materials

This supplement material includes the proofs for Theorem 2.

Because the confusion matrix W is symmetric, it can be decomposed as $W = P\Lambda P^\top$, where $P = (v_1, v_2, \dots, v_n)$ is an orthogonal matrix, i.e., $P^\top P = PP^\top = I$, and $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_n\}$ is a diagonal matrix with diagonal entries being the eigenvalues of W in the nonincreasing order. Then applying the decomposition to the iteration (from W_t and W_{t-1} to W_{t+1})

$$X_{t+1} = 2X_t W - X_{t-1} W - \gamma G(X_t; \xi_t) W + \gamma G(X_{t-1}; \xi_{t-1}) W$$

gives

$$X_{t+1} = 2X_t P \Lambda P^\top - X_{t-1} P \Lambda P^\top - \gamma G(X_t; \xi_t) P \Lambda P^\top + \gamma G(X_{t-1}; \xi_{t-1}) P \Lambda P^\top.$$

Denote $Y_t = X_t P$, $H(X_t; \xi_t) = G(X_t; \xi_t) P$, and use $y_t^{(i)}$ and $h_t^{(i)}$ to indicate the i -th column of Y_t and $H(X_t; \xi_t)$, respectively. Then

$$Y_{t+1} = 2Y_t \Lambda - Y_{t-1} \Lambda - \gamma H(X_t; \xi_t) \Lambda + \gamma H(X_{t-1}; \xi_{t-1}) \Lambda, \quad (6)$$

or in the columns of Y_t and $H(X_t; \xi_t)$,

$$y_{t+1}^{(i)} = \lambda_i (2y_t^{(i)} - y_{t-1}^{(i)} - \gamma h_t^{(i)} + \gamma h_{t-1}^{(i)}). \quad (7)$$

From the properties of W in Assumption 1 and the decomposition, we have $\lambda_1 = 1$ and $v_1 = \frac{1}{\sqrt{n}}(1, 1, \dots, 1)^\top$. Therefore $y_t^{(1)} = \bar{X}_t \sqrt{n}$. For all other eigenvalues $-\frac{1}{3} < \lambda_i < 1$, the equation (7) shows that all $y_t^{(i)}$ would "decay to zero", which explains how the confusion matrix works.

Lemma 4. *Given two non-negative sequences $\{a_t\}_{t=1}^\infty$ and $\{b_t\}_{t=1}^\infty$ that satisfying*

$$a_t = \sum_{s=1}^t \rho^{t-s} b_s, \quad (8)$$

with $\rho \in [0, 1)$, we have

$$S_k := \sum_{t=1}^k a_t \leq \sum_{s=1}^k \frac{b_s}{1-\rho},$$

$$D_k := \sum_{t=1}^k a_t^2 \leq \frac{1}{(1-\rho)^2} \sum_{s=1}^k b_s^2.$$

Proof.

$$S_k := \sum_{t=1}^k a_t = \sum_{t=1}^k \sum_{s=1}^t \rho^{t-s} b_s = \sum_{s=1}^k \sum_{t=s}^k \rho^{t-s} b_s = \sum_{s=1}^k \sum_{t=0}^{k-s} \rho^t b_s \leq \sum_{s=1}^k \frac{b_s}{1-\rho}. \quad (9)$$

$$D_k := \sum_{t=1}^k a_t^2 = \sum_{t=1}^k \sum_{s=1}^t \rho^{t-s} b_s \sum_{r=1}^t \rho^{t-r} b_r = \sum_{t=1}^k \sum_{s=1}^t \sum_{r=1}^t \rho^{2t-s-r} b_s b_r$$

$$\begin{aligned}
&\leq \sum_{t=1}^k \sum_{s=1}^t \sum_{r=1}^t \rho^{2t-s-r} \frac{b_s^2 + b_r^2}{2} = \sum_{t=1}^k \sum_{s=1}^t \sum_{r=1}^t \rho^{2t-s-r} b_s^2 \\
&\leq \frac{1}{1-\rho} \sum_{t=1}^k \sum_{s=1}^t \rho^{t-s} b_s^2 \leq \frac{1}{(1-\rho)^2} \sum_{s=1}^k b_s^2
\end{aligned} \tag{10}$$

where the last inequality holds because of (9). \square

Lemma 5. For any matrix $X_t \in \mathbb{R}^{N \times n}$, we have

$$\begin{aligned}
\sum_{i=2}^n \left\| X_t \mathbf{v}^{(i)} \right\|^2 &\leq \sum_{i=1}^n \left\| X_t \mathbf{v}^{(i)} \right\|^2 = \|X_t\|_F^2 \\
\sum_{i=1}^n \left\| X_t P^\top \mathbf{e}^{(i)} \right\|^2 &= \left\| X_t P^\top \right\|_F^2 = \|X_t\|_F^2
\end{aligned}$$

where $\mathbf{e}^{(i)} \in \mathbb{R}^{n \times 1}$ with the i -th component being 1 and all others being 0.

Proof. From the definition of the Frobenius norm for a matrix, we have

$$\sum_{i=1}^n \left\| X_t \mathbf{v}^{(i)} \right\|^2 = \|X_t P\|_F^2 = \text{Tr} \left(X_t P P^\top X_t^\top \right) = \text{Tr} \left(X_t X_t^\top \right) = \|X_t\|_F^2.$$

Since $\left\| X_t \mathbf{v}^{(1)} \right\|^2 \geq 0$, so

$$\sum_{i=2}^n \left\| X_t \mathbf{v}^{(i)} \right\|^2 \leq \sum_{i=1}^n \left\| X_t \mathbf{v}^{(i)} \right\|^2 = \|X_t\|_F^2.$$

In the same way, we have

$$\sum_{i=1}^n \left\| X_t P^\top \mathbf{e}^{(i)} \right\|^2 = \left\| X_t P^\top \right\|_F^2 = \|X_t\|_F^2.$$

The result is proved. \square

Lemma 6. Given $\rho \in (-\frac{1}{3}, 0) \cup (0, 1)$, for any two sequence $\{a_t\}_{t=0}^\infty$ and $\{b_t\}_{t=0}^\infty$ that satisfy

$$\begin{aligned}
a_0 &= b_0 = 0, \\
a_1 &= b_1, \\
a_{t+1} &= \rho(2a_t - a_{t-1}) + b_t - b_{t-1}, \quad \forall t \geq 1,
\end{aligned}$$

we have

$$a_{t+1} = a_1 \left(\frac{u^{t+1} - v^{t+1}}{u - v} \right) + \sum_{s=1}^t \beta_s \frac{u^{t-s+1} - v^{t-s+1}}{u - v}, \quad \forall t \geq 0,$$

where

$$\beta_s = b_s - b_{s-1}, \quad u = \rho + \sqrt{\rho^2 - \rho}, \quad v = \rho - \sqrt{\rho^2 - \rho}.$$

More specifically, if $0 < \rho < 1$, we have

$$a_{t+1} \sin \theta = a_1 \rho^{t/2} \sin [(t+1)\theta] + \sum_{s=1}^t \beta_s \rho^{(t-s)/2} \sin [(t+1-s)\theta], \quad \forall t \geq 0$$

where

$$\beta_s = b_s - b_{s-1}, \quad \theta = \arccos(\sqrt{\rho}).$$

Proof. When $t = 0$, the results is easy to verify. Next we consider the case $t \geq 1$. Since

$$a_{t+1} = 2\rho a_t - \rho a_{t-1} + \beta_t,$$

We can find

$$u = \rho + \sqrt{\rho^2 - \rho}, \quad v = \rho - \sqrt{\rho^2 - \rho},$$

such that

$$a_{t+1} - ua_t = (a_t - ua_{t-1})v + \beta_t. \tag{11}$$

Note that u and v are complex numbers when $0 < \rho < 1$. That is

$$u = \sqrt{\rho}e^{i\theta}, \quad v = \sqrt{\rho}e^{-i\theta},$$

with $\theta = \arccos(\sqrt{\rho})$.

Recursively applying (11) gives

$$\begin{aligned} a_{t+1} - ua_t &= (a_t - ua_{t-1})v + \beta_t = (a_{t-1} - ua_{t-2})v^2 + \beta_{t-1}v + \beta_t \\ &= (a_1 - ua_0)v^t + \sum_{s=1}^t \beta_s v^{t-s} \\ &= a_1 v^t + \sum_{s=1}^t \beta_s v^{t-s}. \quad (\text{due to } a_0 = 0) \end{aligned}$$

Dividing both sides by u^{t+1} , we obtain

$$\begin{aligned} \frac{a_{t+1}}{u^{t+1}} &= \frac{a_t}{u^t} + u^{-(t+1)} \left(a_1 v^t + \sum_{s=1}^t \beta_s v^{t-s} \right) \\ &= \frac{a_{t-1}}{u^{t-1}} + u^{-t} \left(a_1 v^{t-1} + \sum_{s=1}^{t-1} \beta_s v^{t-1-s} \right) + u^{-(t+1)} \left(a_1 v^t + \sum_{s=1}^t \beta_s v^{t-s} \right) \\ &= \frac{a_1}{u} + \sum_{k=1}^t u^{-k-1} \left(a_1 v^k + \sum_{s=1}^k \beta_s v^{k-s} \right) \end{aligned}$$

Then we multiply both sides by u^{t+1} and have

$$\begin{aligned} a_{t+1} &= a_1 u^t + \sum_{k=1}^t u^{t-k} \left(a_1 v^k + \sum_{s=1}^k \beta_s v^{k-s} \right) \\ &= a_1 u^t \left(1 + \sum_{k=1}^t \left(\frac{v}{u} \right)^k \right) + u^t \sum_{k=1}^t \sum_{s=1}^k \beta_s v^{-s} \left(\frac{v}{u} \right)^k \\ &= a_1 u^t \sum_{k=0}^t \left(\frac{v}{u} \right)^k + u^t \sum_{s=1}^t \sum_{k=s}^t \beta_s v^{-s} \left(\frac{v}{u} \right)^k \quad (\text{due to } \sum_{k=1}^t \sum_{s=1}^k a_s b_k = \sum_{s=1}^t \sum_{k=s}^t a_s b_k) \end{aligned}$$

$$\begin{aligned}
&= a_1 u^t \left(\frac{1 - \left(\frac{v}{u}\right)^{t+1}}{1 - \frac{v}{u}} \right) + u^t \sum_{s=1}^t \beta_s v^{-s} \left(\frac{v}{u}\right)^s \frac{1 - \left(\frac{v}{u}\right)^{t-s+1}}{1 - \frac{v}{u}} \\
&= a_1 \left(\frac{u^{t+1} - v^{t+1}}{u - v} \right) + \sum_{s=1}^t \beta_s \frac{u^{t-s+1} - v^{t-s+1}}{u - v}.
\end{aligned}$$

When $\rho \in (0, 1)$, since $u = \sqrt{\rho}e^{i\theta}$ and $v = \sqrt{\rho}e^{-i\theta}$, we have

$$a_{t+1} = a_1 \rho^{t/2} \frac{\sin[(t+1)\theta]}{\sin\theta} + \sum_{s=1}^t \beta_s \rho^{(t-s)/2} \frac{\sin[(t-s+1)\theta]}{\sin\theta}.$$

The result is proved. □

Lemma 7. *Under Assumption 1, we have*

$$\begin{aligned}
&\left(1 - 24C_2\gamma^2L^2\right) \sum_{i=1}^n \sum_{t=0}^T \left\| \bar{X}_t - \mathbf{x}_t^{(i)} \right\|^2 \\
&\leq 2C_1 \|\mathbf{X}_1\|_F^2 + 12C_2\gamma^2n\sigma^2T + 6C_2\gamma^4L^2\sigma^2T + 6C_2\gamma^4L^2n \sum_{t=1}^{T-1} \left\| \nabla f(\mathbf{X}_t) \right\|^2,
\end{aligned}$$

where $\gamma, L, \sigma, \theta, C_1$ and C_2 are defined in Theorem 2.

Proof. To estimate the difference of the local models and the global mean model, we have

$$\begin{aligned}
\sum_{i=1}^n \left\| \bar{X}_t - \mathbf{x}_t^{(i)} \right\|^2 &= \sum_{i=1}^n \left\| X_t \mathbf{e}^{(i)} - X_t \frac{\mathbf{1}_n}{n} \right\|^2 = \left\| X_t - X_t \frac{\mathbf{1}_n \mathbf{1}_n^\top}{n} \right\|_F^2 = \left\| X_t P P^\top - X_t \mathbf{v}_1 \mathbf{v}_1^\top \right\|_F^2 \\
&= \left\| X_t P \begin{pmatrix} 0, & 0, & 0, & \dots, & 0 \\ 0, & 1, & 0, & \dots, & 0 \\ 0, & 0, & 1, & \dots, & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0, & 0, & 0, & \dots, & 1 \end{pmatrix} \right\|_F^2 \\
&= \sum_{i=2}^n \left\| \mathbf{y}_t^{(i)} \right\|^2,
\end{aligned} \tag{12}$$

where $\mathbf{y}_t^{(i)}$ is the i -th column of $X_t P$. Note that we have, from (7),

$$\mathbf{y}_{t+1}^{(i)} = \lambda_i (2\mathbf{y}_t^{(i)} - \mathbf{y}_{t-1}^{(i)} - \gamma h_t^{(i)} + \gamma h_{t-1}^{(i)}) = \lambda_i (2\mathbf{y}_t^{(i)} - \mathbf{y}_{t-1}^{(i)}) + \lambda_i \beta_t^{(i)},$$

where $\beta_t^{(i)} = -\gamma h_t^{(i)} + \gamma h_{t-1}^{(i)}$. For all $\mathbf{y}^{(i)}$ that corresponding to $-\frac{1}{3} < \lambda_i < 0$, Lemma 6 shows

$$\mathbf{y}_{t+1}^{(i)} = \mathbf{y}_1^{(i)} \left(\frac{u_i^{t+1} - v_i^{t+1}}{u_i - v_i} \right) + \lambda_i \sum_{s=1}^t \beta_s^{(i)} \frac{u_i^{t-s+1} - v_i^{t-s+1}}{u_i - v_i},$$

where $u_i = \lambda_i + \sqrt{\lambda_i^2 - \lambda_i}$ and $v_i = \lambda_i - \sqrt{\lambda_i^2 - \lambda_i}$. Therefore, we have

$$\left\| \mathbf{y}_{t+1}^{(i)} \right\|^2 \leq 2 \left\| \mathbf{y}_1^{(i)} \right\|^2 \left(\frac{u_i^{t+1} - v_i^{t+1}}{u_i - v_i} \right)^2 + 2\lambda_i^2 \left(\sum_{s=1}^t \left\| \beta_s^{(i)} \right\| \left| \frac{u_i^{t-s+1} - v_i^{t-s+1}}{u_i - v_i} \right| \right)^2. \tag{13}$$

For $\left| \frac{u_i^{n+1} - v_i^{n+1}}{u_i - v_i} \right|$, we have

$$\left| \frac{u_i^{n+1} - v_i^{n+1}}{u_i - v_i} \right| \leq |v_i|^n \left| \frac{u_i \left(\frac{u_i}{v_i} \right)^n - v_i}{u_i - v_i} \right| \leq |v_i|^n \quad (\text{due to } |u_i| < |v_i|).$$

Using (13), we obtain

$$\|y_{t+1}^{(i)}\|^2 \leq 2 \|y_1^{(i)}\| |v_i|^{2t} + 2\lambda_i^2 \left(\sum_{s=1}^t \|\beta_s^{(i)}\| |v_i|^{t-s} \right)^2.$$

Summing from $t = 0$ to $t = T - 1$ gives

$$\sum_{t=0}^{T-1} \|y_{t+1}^{(i)}\|^2 = \sum_{t=1}^T \|y_t^{(i)}\|^2 \leq 2 \|y_1^{(i)}\| \sum_{t=0}^{T-1} |v_i|^{2t} + 2\lambda_i^2 \sum_{t=1}^{T-1} \left(\sum_{s=1}^t \|\beta_s^{(i)}\| |v_i|^{t-s} \right)^2.$$

Denote $a_t = \sum_{s=1}^t \|\beta_s^{(i)}\| |v_i|^{t-s}$, which has the same structure as the sequence in Lemma 4. Therefore, when $\lambda_i < 0$, we have

$$\begin{aligned} \sum_{t=1}^T \|y_t^{(i)}\|^2 &\leq \frac{2 \|y_1^{(i)}\|}{1 - |v_i|^2} + \frac{2\lambda_i^2}{(1 - |v_i|)^2} \sum_{t=1}^{T-1} \|\beta_t^{(i)}\|^2 \\ &\leq \frac{2 \|y_1^{(i)}\|}{1 - |v|^2} + \frac{2\lambda_n^2}{(1 - |v|)^2} \sum_{t=1}^{T-1} \|\beta_t^{(i)}\|^2, \end{aligned} \quad (14)$$

where $v = \lambda_n - \sqrt{\lambda_n^2 - \lambda_n}$.

For all $y^{(i)}$ that satisfies $0 \leq \lambda_i < 1$, from (7) and Lemma 6, we have

$$y_{t+1}^{(i)} \sin \theta_i = y_1^{(i)} \lambda_i^{t/2} \sin [(t+1)\theta_i] + \lambda_i \sum_{s=1}^t \beta_s^{(i)} \lambda_i^{(t-s)/2} \sin [(t+1-s)\theta_i],$$

where $\beta_s^{(i)} = -\gamma h_s^{(i)} + \gamma h_{s-1}^{(i)}$ and $\theta_i = \arccos(\sqrt{\lambda_i})$.

Then

$$\begin{aligned} \|y_{t+1}^{(i)}\|^2 \sin^2 \theta_i &\leq 2 \|y_1^{(i)}\|^2 \lambda_i^t \sin^2 [(t+1)\theta_i] + 2\lambda_i^2 \left(\sum_{s=1}^t \|\beta_s^{(i)}\| \sin [(t+1-s)\theta_i] \|\lambda_i^{(t-s)/2}\| \right)^2 \\ &\leq 2 \|y_1^{(i)}\|^2 \lambda_i^t + 2\lambda_i^2 \left(\sum_{s=1}^t \|\beta_s^{(i)}\| \|\lambda_i^{(t-s)/2}\| \right)^2, \end{aligned}$$

Summing from $t = 0$ to $T - 1$ gives

$$\sum_{t=0}^{T-1} \|y_{t+1}^{(i)}\|^2 \sin^2 \theta_i = \sum_{t=1}^T \|y_t^{(i)}\|^2 \sin^2 \theta_i \leq 2 \|y_1^{(i)}\|^2 \sum_{t=0}^{T-1} \lambda_i^t + 2\lambda_i^2 \sum_{t=1}^{T-1} \left(\sum_{s=1}^t \|\beta_s^{(i)}\| \|\lambda_i^{(t-s)/2}\| \right)^2$$

From Lemma 4, $\sum_{s=1}^t \|\beta_s^{(i)}\| \lambda_i^{(t-s)/2}$ has the same structure as the sequence in Lemma 4, so we have

$$\sum_{t=1}^T \left\| y_t^{(i)} \right\|^2 \sin^2 \theta_i \leq \frac{2 \left\| y_1^{(i)} \right\|^2}{1 - \lambda_i} + \frac{2\lambda_i^2}{(1 - \sqrt{\lambda_i})^2} \sum_{t=1}^{T-1} \left\| \beta_t^{(i)} \right\|^2.$$

Then $\sin^2 \theta_i = 1 - \lambda_i$ gives

$$\begin{aligned} \sum_{t=1}^T \left\| y_t^{(i)} \right\|^2 &\leq \frac{2 \left\| y_1^{(i)} \right\|^2}{(1 - \lambda_i)^2} + \frac{2\lambda_i^2}{(1 - \sqrt{\lambda_i})^2 (1 - \lambda_i)} \sum_{t=1}^{T-1} \left\| \beta_t^{(i)} \right\|^2 \\ &\leq \frac{2 \left\| y_1^{(i)} \right\|^2}{(1 - \lambda)^2} + \frac{2\lambda^2}{(1 - \sqrt{\lambda})^2 (1 - \lambda)} \sum_{t=1}^{T-1} \left\| \beta_t^{(i)} \right\|^2. \end{aligned} \quad (15)$$

Denote $C_1 = \max \left\{ \frac{1}{1 - |\nu|^2}, \frac{1}{(1 - \lambda)^2} \right\}$ and $C_2 = \max \left\{ \frac{\lambda_n^2}{(1 - |\nu|^2)}, \frac{\lambda^2}{(1 - \sqrt{\lambda})^2 (1 - \lambda)} \right\}$. From (14) and (15), we have

$$\sum_{t=1}^T \left\| y_t^{(i)} \right\|^2 \leq 2C_1 \left\| y_1^{(i)} \right\|^2 + 2C_2 \sum_{t=1}^{T-1} \left\| \beta_t^{(i)} \right\|^2. \quad (16)$$

We next bound $\beta_t^{(i)}$

$$\begin{aligned} &\mathbb{E} \sum_{i=2}^n \left\| \beta_t^{(i)} \right\|^2 \\ &= \sum_{i=2}^n \gamma^2 \mathbb{E} \left\| h_t^{(i)} - h_{t-1}^{(i)} \right\|^2 \\ &= \gamma^2 \sum_{i=2}^n \mathbb{E} \left\| G(X_t; \xi_t) P e^{(i)} - G(X_{t-1}; \xi_{t-1}) P e^{(i)} \right\|^2 \\ &\leq \gamma^2 \sum_{i=1}^n \mathbb{E} \left\| G(X_t; \xi_t) P e^{(i)} - G(X_{t-1}; \xi_{t-1}) P e^{(i)} \right\|^2 \\ &= \gamma^2 \mathbb{E} \left\| G(X_t; \xi_t) P - G(X_{t-1}; \xi_{t-1}) P \right\|_F^2 \\ &= \gamma^2 \mathbb{E} \left\| G(X_t; \xi_t) - G(X_{t-1}; \xi_{t-1}) \right\|_F^2 \quad (\text{due to Lemma 5}) \\ &= \gamma^2 \sum_{i=1}^n \mathbb{E} \left\| \nabla F_i(\mathbf{x}_t^{(i)}; \xi_t^{(i)}) - \nabla F_i(\mathbf{x}_{t-1}^{(i)}; \xi_{t-1}^{(i)}) \right\|^2 \\ &= \gamma^2 \sum_{i=1}^n \mathbb{E} \left\| \left(\nabla F_i(\mathbf{x}_t^{(i)}; \xi_t^{(i)}) - \nabla f_i(\mathbf{x}_t^{(i)}) \right) - \left(F_i(\mathbf{x}_{t-1}^{(i)}; \xi_{t-1}^{(i)}) - \nabla f_i(\mathbf{x}_{t-1}^{(i)}) \right) + \left(\nabla f_i(\mathbf{x}_t^{(i)}) - \nabla f_i(\mathbf{x}_{t-1}^{(i)}) \right) \right\|^2 \\ &= 3\gamma^2 \sum_{i=1}^n \mathbb{E} \left\| \nabla F_i(\mathbf{x}_t^{(i)}; \xi_t^{(i)}) - \nabla f_i(\mathbf{x}_t^{(i)}) \right\|^2 + 3\gamma^2 \sum_{i=1}^n \left\| F_i(\mathbf{x}_{t-1}^{(i)}; \xi_{t-1}^{(i)}) - \nabla f_i(\mathbf{x}_{t-1}^{(i)}) \right\|^2 \\ &\quad + 3\gamma^2 \sum_{i=1}^n \left\| \nabla f_i(\mathbf{x}_t^{(i)}) - \nabla f_i(\mathbf{x}_{t-1}^{(i)}) \right\|^2 \\ &\leq 6\gamma^2 n \sigma^2 + 3\gamma^2 \sum_{i=1}^n \mathbb{E} \left\| \nabla f_i(\mathbf{x}_t^{(i)}) - \nabla f_i(\mathbf{x}_{t-1}^{(i)}) \right\|^2 \\ &\leq 6\gamma^2 n \sigma^2 + 3\gamma^2 \sum_{i=1}^n L^2 \mathbb{E} \left\| \mathbf{x}_t^{(i)} - \mathbf{x}_{t-1}^{(i)} \right\|^2 \end{aligned}$$

$$\begin{aligned}
&= 6\gamma^2 n\sigma^2 + 3\gamma^2 L^2 \sum_{i=1}^n \mathbb{E} \left\| Y_t P^\top \mathbf{e}^{(i)} - Y_{t-1} P^\top \mathbf{e}^{(i)} \right\|^2 \\
&= 6\gamma^2 n\sigma^2 + 3\gamma^2 L^2 \mathbb{E} \left\| Y_t P^\top - Y_{t-1} P^\top \right\|_F^2 \\
&= 6\gamma^2 n\sigma^2 + 3\gamma^2 L^2 \mathbb{E} \|Y_t - Y_{t-1}\|_F^2 \quad (\text{due to Lemma 5}) \\
&= 6\gamma^2 n\sigma^2 + 3\gamma^2 L^2 \sum_{i=1}^n \mathbb{E} \left\| \mathbf{y}_t^{(i)} - \mathbf{y}_{t-1}^{(i)} \right\|^2. \tag{17}
\end{aligned}$$

Combing (16) and (17), we have

$$\begin{aligned}
\sum_{i=2}^n \sum_{t=1}^T \|\mathbf{y}_t^{(i)}\|^2 &\leq 2C_1 \|Y_1\|_F^2 + 2C_2 \sum_{i=2}^n \sum_{t=1}^{T-1} \|\beta_t^{(i)}\|^2 \\
&\leq 2C_1 \|Y_1\|_F^2 + 2C_2 \sum_{t=1}^{T-1} \left(6\gamma^2 n\sigma^2 + 3\gamma^2 L^2 \sum_{i=1}^n \mathbb{E} \left\| \mathbf{y}_t^{(i)} - \mathbf{y}_{t-1}^{(i)} \right\|^2 \right) \\
&\leq 2C_1 \|Y_1\|_F^2 + 12C_2 \gamma^2 n\sigma^2 T + 6C_2 \gamma^2 L^2 \sum_{i=1}^n \sum_{t=1}^{T-1} \mathbb{E} \left\| \mathbf{y}_t^{(i)} - \mathbf{y}_{t-1}^{(i)} \right\|^2. \tag{18}
\end{aligned}$$

The next step is to bound $\mathbb{E} \|\mathbf{y}_t^{(1)} - \mathbf{y}_{t-1}^{(1)}\|^2$. Because

$$\mathbf{y}_t^{(1)} = X_t P \mathbf{e}^{(1)} = X_t \mathbf{v}_1 = X_t \frac{1}{\sqrt{n}} \mathbf{1}_n = \bar{X}_t \sqrt{n},$$

what we need to bound is $\mathbb{E} \|\bar{X}_{t+1} - \bar{X}_t\|^2$. From (4), we have $\bar{X}_{t+1} = \bar{X}_t - \gamma \bar{G}_t$. Therefore

$$\mathbb{E} \|\bar{X}_{t+1} - \bar{X}_t\|^2 = \gamma^2 \mathbb{E} \|\bar{G}_t\|^2 = \gamma^2 \mathbb{E} \|\bar{G}_t - \bar{\nabla} f(X_t)\|^2 + \gamma^2 \|\bar{\nabla} f(X_t)\|^2 \leq \frac{\gamma^2 \sigma^2}{n} + \gamma^2 \|\bar{\nabla} f(X_t)\|^2,$$

and we have the follow bound for $\mathbb{E} \|\mathbf{y}_t^{(1)} - \mathbf{y}_{t-1}^{(1)}\|^2$:

$$\mathbb{E} \left\| \mathbf{y}_{t+1}^{(1)} - \mathbf{y}_t^{(1)} \right\|^2 \leq \gamma^2 \sigma^2 + n\gamma^2 \|\bar{\nabla} f(X_t)\|^2. \tag{19}$$

Combing (18) and (19) we get

$$\begin{aligned}
\sum_{i=2}^n \sum_{t=1}^T \|\mathbf{y}_t^{(i)}\|^2 &\leq 2C_1 \|Y_1\|_F^2 + 12C_2 \gamma^2 n\sigma^2 T + 6C_2 \gamma^4 L^2 \sigma^2 T + 6C_2 \gamma^4 L^2 n \sum_{t=1}^{T-1} \|\bar{\nabla} f(X_t)\|^2 \\
&\quad + 6C_2 \gamma^2 L^2 \sum_{i=2}^n \sum_{t=1}^{T-1} \mathbb{E} \left\| \mathbf{y}_t^{(i)} - \mathbf{y}_{t-1}^{(i)} \right\|^2 \\
&\leq 2C_1 \|Y_1\|_F^2 + 12C_2 \gamma^2 n\sigma^2 T + 6C_2 \gamma^4 L^2 \sigma^2 T + 6C_2 \gamma^4 L^2 n \sum_{t=1}^{T-1} \|\bar{\nabla} f(X_t)\|^2 \\
&\quad + 6C_2 \gamma^2 L^2 \sum_{i=2}^n \sum_{t=1}^{T-1} 2\mathbb{E} \left(\left\| \mathbf{y}_t^{(i)} \right\|^2 + \left\| \mathbf{y}_{t-1}^{(i)} \right\|^2 \right) \\
&\leq 2C_1 \|Y_1\|_F^2 + 12C_2 \gamma^2 n\sigma^2 T + 6C_2 \gamma^4 L^2 \sigma^2 T + 6C_2 \gamma^4 L^2 n \sum_{t=1}^{T-1} \|\bar{\nabla} f(X_t)\|^2 \\
&\quad + 6C_2 \gamma^2 L^2 \sum_{i=2}^n \sum_{t=1}^{T-1} 2\mathbb{E} \left(\left\| \mathbf{y}_t^{(i)} \right\|^2 + \left\| \mathbf{y}_{t-1}^{(i)} \right\|^2 \right) \quad (\text{due to } \mathbf{y}_0^{(i)} = \mathbf{0})
\end{aligned}$$

$$\begin{aligned}
&\leq 2C_1 \|Y_1\|_F^2 + 12C_2 \gamma^2 n \sigma^2 T + 6C_2 \gamma^4 L^2 \sigma^2 T + 6C_2 \gamma^4 L^2 n \sum_{t=1}^{T-1} \|\overline{\nabla f}(X_t)\|^2 \\
&\quad + 24C_2 \gamma^2 L^2 \sum_{i=2}^n \sum_{t=1}^{T-1} \mathbb{E} \|y_t^{(i)}\|^2, \\
(1 - 24C_2 \gamma^2 L^2) \sum_{i=2}^n \sum_{t=1}^T \|y_t^{(i)}\|^2 &\leq 2C_1 \|Y_1\|_F^2 + 12C_2 \gamma^2 n \sigma^2 T + 6C_2 \gamma^4 L^2 \sigma^2 T + 6C_2 \gamma^4 L^2 n \sum_{t=1}^{T-1} \|\overline{\nabla f}(X_t)\|^2.
\end{aligned}$$

Together with (12) and $X_0 = 0$, we have

$$\begin{aligned}
(1 - 24C_2 \gamma^2 L^2) \sum_{i=1}^n \sum_{t=1}^T \|\bar{X}_t - \mathbf{x}_t^{(i)}\|^2 &\leq 2C_1 \|Y_1\|_F^2 + 12C_2 \gamma^2 n \sigma^2 T + 6C_2 \gamma^4 L^2 \sigma^2 T + 6C_2 \gamma^4 L^2 n \sum_{t=1}^{T-1} \|\overline{\nabla f}(X_t)\|^2 \\
(\text{due to } \|X_1\|_F = \|Y_1\|_F) &\leq 2C_1 \|X_1\|_F^2 + 12C_2 \gamma^2 n \sigma^2 T + 6C_2 \gamma^4 L^2 \sigma^2 T + 6C_2 \gamma^4 L^2 n \sum_{t=1}^{T-1} \|\overline{\nabla f}(X_t)\|^2.
\end{aligned}$$

Actually, when $\lambda_n \leq -\frac{1}{3}$, we have $|v_n| \geq 1$, then $\|y_t^{(n)}\|^2 \propto t$ and

$$\frac{1}{T} \sum_{i=1}^n \sum_{t=1}^T \|\bar{X}_t - \mathbf{x}_t^{(i)}\|^2 \leq T.$$

The algorithm would fail to converge in this situation, and this is why $-1/3$ is the infimum of λ_n . \square

Lemma 8. *Following the Assumption 1, we have*

$$\mathbb{E}f(\bar{X}_{t+1}) \leq \mathbb{E}f(\bar{X}_t) - \frac{\gamma_t}{2} \mathbb{E}\|\nabla f(\bar{X}_t)\|^2 - \left(\frac{\gamma_t}{2} - \frac{L\gamma_t^2}{2}\right) \mathbb{E}\|\overline{\nabla f}(X_t)\|^2 + \frac{\gamma_t}{2} \mathbb{E}\|\nabla f(\bar{X}_t) - \overline{\nabla f}(X_t)\|^2 + \frac{L\gamma_t^2}{2n} \sigma^2.$$

Proof. From (4), we have

$$\bar{X}_{t+1} = \bar{X}_t - \gamma_t \bar{G}(X_t; \xi_t).$$

From item 1 of Assumption 1, we know that f has a L -Lipschitz continuous gradient. So, we have

$$\begin{aligned}
\mathbb{E}f(\bar{X}_{t+1}) &\leq \mathbb{E}f(\bar{X}_t) + \mathbb{E}\langle \nabla f(\bar{X}_t), -\gamma_t \bar{G}(X_t; \xi_t) \rangle + \frac{L}{2} \mathbb{E}\|-\gamma_t \bar{G}(X_t; \xi_t)\|^2 \\
&= \mathbb{E}f(\bar{X}_t) + \mathbb{E}\langle \nabla f(\bar{X}_t), -\gamma_t \mathbb{E}_{\xi_t} \bar{G}(X_t; \xi_t) \rangle + \frac{L\gamma_t^2}{2} \mathbb{E}\|\bar{G}(X_t; \xi_t)\|^2 \\
&= \mathbb{E}f(\bar{X}_t) - \gamma_t \mathbb{E}\langle \nabla f(\bar{X}_t), \overline{\nabla f}(X_t) \rangle + \frac{L\gamma_t^2}{2} \mathbb{E}\|(\bar{G}(X_t; \xi_t) - \overline{\nabla f}(X_t)) + \overline{\nabla f}(X_t)\|^2 \\
&= \mathbb{E}f(\bar{X}_t) - \gamma_t \mathbb{E}\langle \nabla f(\bar{X}_t), \overline{\nabla f}(X_t) \rangle + \frac{L\gamma_t^2}{2} \mathbb{E}\|\bar{G}(X_t; \xi_t) - \overline{\nabla f}(X_t)\|^2 + \frac{L\gamma_t^2}{2} \mathbb{E}\|\overline{\nabla f}(X_t)\|^2 \\
&\quad + L\gamma_t^2 \mathbb{E}\langle \mathbb{E}_{\xi_t} \bar{G}(X_t; \xi_t) - \overline{\nabla f}(X_t), \overline{\nabla f}(X_t) \rangle \\
&= \mathbb{E}f(\bar{X}_t) - \gamma_t \mathbb{E}\langle \nabla f(\bar{X}_t), \overline{\nabla f}(X_t) \rangle + \frac{L\gamma_t^2}{2} \mathbb{E}\|\bar{G}(X_t; \xi_t) - \overline{\nabla f}(X_t)\|^2 + \frac{L\gamma_t^2}{2} \mathbb{E}\|\overline{\nabla f}(X_t)\|^2 \\
&= \mathbb{E}f(\bar{X}_t) - \gamma_t \mathbb{E}\langle \nabla f(\bar{X}_t), \overline{\nabla f}(X_t) \rangle + \frac{L\gamma_t^2}{2n^2} \mathbb{E}\left\| \sum_{i=1}^n (\nabla F_i(x_t^{(i)}; \xi_t^{(i)}) - \nabla f_i(x_t^{(i)})) \right\|^2 \\
&\quad + \frac{L\gamma_t^2}{2} \mathbb{E}\|\overline{\nabla f}(X_t)\|^2
\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}f(\bar{X}_t) - \gamma_t \mathbb{E} \langle \nabla f(\bar{X}_t), \bar{\nabla} f(X_t) \rangle + \frac{L\gamma_t^2}{2n^2} \sum_{i=1}^n \mathbb{E} \left\| \nabla F_i(x_t^{(i)}; \xi_t^{(i)}) - \nabla f_i(x_t^{(i)}) \right\|^2 \\
&\quad + \sum_{i \neq i'}^n \mathbb{E} \left\langle \mathbb{E}_{\xi_t} \nabla F_i(x_t^{(i)}; \xi_t^{(i)}) - \nabla f_i(x_t^{(i)}), \nabla \mathbb{E}_{\xi_t} F_{i'}(x_t^{(i')}; \xi_t^{(i')}) - \nabla f_{i'}(x_t^{(i')}) \right\rangle + \frac{L\gamma_t^2}{2} \mathbb{E} \|\bar{\nabla} f(X_t)\|^2 \\
&\leq \mathbb{E}f(\bar{X}_t) - \gamma_t \mathbb{E} \langle \nabla f(\bar{X}_t), \bar{\nabla} f(X_t) \rangle + \frac{L\gamma_t^2}{2n} \sigma^2 + \frac{L\gamma_t^2}{2} \mathbb{E} \|\bar{\nabla} f(X_t)\|^2 \\
&= \mathbb{E}f(\bar{X}_t) - \frac{\gamma_t}{2} \mathbb{E} \|\nabla f(\bar{X}_t)\|^2 - \frac{\gamma_t}{2} \mathbb{E} \|\bar{\nabla} f(X_t)\|^2 + \frac{\gamma_t}{2} \mathbb{E} \|\nabla f(\bar{X}_t) - \bar{\nabla} f(X_t)\|^2 + \frac{L\gamma_t^2}{2} \mathbb{E} \|\bar{\nabla} f(X_t)\|^2 \\
&\quad + \frac{L\gamma_t^2}{2n} \sigma^2 \quad (\text{due to } 2\langle \mathbf{a}, \mathbf{b} \rangle = \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 - \|\mathbf{a} - \mathbf{b}\|^2) \\
&= \mathbb{E}f(\bar{X}_t) - \frac{\gamma_t}{2} \mathbb{E} \|\nabla f(\bar{X}_t)\|^2 - \left(\frac{\gamma_t}{2} - \frac{L\gamma_t^2}{2} \right) \mathbb{E} \|\bar{\nabla} f(X_t)\|^2 + \frac{\gamma_t}{2} \mathbb{E} \|\nabla f(\bar{X}_t) - \bar{\nabla} f(X_t)\|^2 + \frac{L\gamma_t^2}{2n} \sigma^2,
\end{aligned} \tag{20}$$

which completes the proof. \square

Proof to Theorem 2

Proof. We first estimate the upper bound for $\mathbb{E} \|\nabla f(\bar{X}_t) - \bar{\nabla} f(X_t)\|^2$:

$$\begin{aligned}
\mathbb{E} \|\nabla f(\bar{X}_t) - \bar{\nabla} f(X_t)\|^2 &= \frac{1}{n^2} \mathbb{E} \left\| \sum_{i=1}^n \left(\nabla f_i(\bar{X}_t) - \nabla f_i(x_t^{(i)}) \right) \right\|^2 \\
&\leq \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left\| \nabla f_i(\bar{X}_t) - \nabla f_i(x_t^{(i)}) \right\|^2 \\
&\leq \frac{L^2}{n} \mathbb{E} \sum_{i=1}^n \|\bar{X}_t - x_t^{(i)}\|^2.
\end{aligned} \tag{21}$$

Combining (20) in Lemma 8 and (21) yields

$$\begin{aligned}
&\frac{\gamma_t}{2} \mathbb{E} \|\nabla f(\bar{X}_t)\|^2 + \left(\frac{\gamma_t}{2} - \frac{L\gamma_t^2}{2} \right) \mathbb{E} \|\bar{\nabla} f(X_t)\|^2 \\
&\leq \mathbb{E}f(\bar{X}_t) - \mathbb{E}f(\bar{X}_{t+1}) + \frac{\gamma_t}{2} \mathbb{E} \|\nabla f(\bar{X}_t) - \bar{\nabla} f(X_t)\|^2 + \frac{L\gamma_t^2}{2n} \sigma^2 \\
&\leq \mathbb{E}f(\bar{X}_t) - \mathbb{E}f(\bar{X}_{t+1}) + \frac{L^2\gamma_t}{2n} \sum_{i=1}^n \|\bar{X}_t - x_t^{(i)}\|^2 + \frac{L\gamma_t^2}{2n} \sigma^2.
\end{aligned}$$

Setting $\gamma_t = \gamma$, we obtain

$$\mathbb{E} \|\nabla f(\bar{X}_t)\|^2 + (1 - L\gamma) \mathbb{E} \|\bar{\nabla} f(X_t)\|^2 \leq \frac{2}{\gamma} (\mathbb{E}f(\bar{X}_t) - f^* - (\mathbb{E}f(\bar{X}_{t+1}) - f^*)) + \frac{L^2}{n} \sum_{i=1}^n \|\bar{X}_t - x_t^{(i)}\|^2 + \frac{L\gamma}{n} \sigma^2. \tag{22}$$

From Lemma 7, we have

$$\left(1 - 24C_2\gamma^2L^2\right) \sum_{i=1}^n \sum_{t=0}^T \|\bar{X}_t - x_t^{(i)}\|^2 \leq 2C_1 \|X_1\|_F^2 + 12C_2\gamma^2n\sigma^2T + 6C_2\gamma^4L^2\sigma^2T$$

$$+ 6C_2\gamma^4L^2n \sum_{t=1}^{T-1} \left\| \overline{\nabla f}(X_t) \right\|^2,$$

If γ is not too large that satisfies $1 - 24C_2\gamma^2L^2 > 0$, then denote $C_3 = 1 - 24C_2\gamma^2L^2$, we would have

$$\sum_{i=1}^n \sum_{t=0}^T \left\| \overline{X}_t - x_t^{(i)} \right\|^2 \leq \frac{2C_1}{C_3} \|X_1\|_F^2 + \frac{12C_2\gamma^2n\sigma^2T}{C_3} + \frac{6C_2\gamma^4L^2\sigma^2T}{C_3} + \frac{6C_2\gamma^4L^2n}{C_3} \sum_{t=1}^{T-1} \left\| \overline{\nabla f}(X_t) \right\|^2. \quad (23)$$

Summarizing both sides of (22) and applying (23) yields

$$\begin{aligned} & \sum_{t=0}^{T-1} \left(\mathbb{E} \|\nabla f(\overline{X}_t)\|^2 + (1 - L\gamma) \mathbb{E} \|\overline{\nabla f}(X_t)\|^2 \right) \\ & \leq \frac{2\mathbb{E}f(\overline{X}_0) - 2f^*}{\gamma} + \frac{L^2}{n} \sum_{t=0}^T \sum_{i=1}^n \mathbb{E} \left\| \overline{X}_t - x_t^{(i)} \right\|^2 + \frac{LT\gamma\sigma^2}{n} \\ & \leq \frac{2(f(0) - f^*)}{\gamma} + \frac{LT\gamma\sigma^2}{n} + \frac{2L^2C_1}{nC_3} \|X_1\|_F^2 + \frac{12L^2C_2\gamma^2\sigma^2T}{C_3} + \frac{6L^2C_2\gamma^4L^2\sigma^2T}{nC_3} \\ & \quad + \frac{6L^2C_2\gamma^4L^2}{C_3} \sum_{t=1}^{T-1} \left\| \overline{\nabla f}(X_t) \right\|^2. \end{aligned}$$

It implies

$$\begin{aligned} & \sum_{t=0}^{T-1} \left(\mathbb{E} \|\nabla f(\overline{X}_t)\|^2 + \left(1 - L\gamma - \frac{6L^2C_2\gamma^4L^2}{C_3} \right) \mathbb{E} \|\overline{\nabla f}(X_t)\|^2 \right) \\ & \leq \frac{2(f(0) - f^*)}{\gamma} + \frac{LT\gamma\sigma^2}{n} + \frac{2L^2C_1}{nC_3} \|X_1\|_F^2 + \frac{12L^2C_2\gamma^2n\sigma^2T}{nC_3} + \frac{6L^2C_2\gamma^4L^2\sigma^2T}{nC_3} \\ & = \frac{2(f(0) - f^*)}{\gamma} + \frac{LT\gamma\sigma^2}{n} + \frac{2L^2C_1\gamma^2}{nC_3} \|G(0; \xi_0)\|_F^2 + \frac{12L^2C_2\gamma^2\sigma^2T}{C_3} + \frac{6L^2C_2\gamma^4L^2\sigma^2T}{nC_3}. \end{aligned} \quad (24)$$

However, $\|G(0; \xi_0)\|_F^2$ can be expanded as:

$$\begin{aligned} \|G(0, \xi_0)\|_F^2 &= \sum_{i=1}^n \left\| (\nabla F_i(\mathbf{0}, \xi_1) - \nabla f_i(\mathbf{0})) + (\nabla f_i(\mathbf{0}) - \nabla f(\mathbf{0})) + \nabla f(\mathbf{0}) \right\|^2 \\ &\leq 3n\sigma^2 + 3n\zeta_0^2 + 3n\|\nabla f(\mathbf{0})\|^2, \end{aligned} \quad (25)$$

where $\zeta_0 = \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{0}) - \nabla f(\mathbf{0})\|^2$ indicates the difference between different workers' dataset at the start point. Combining (24) and (25), then we have

$$\begin{aligned} & \sum_{t=0}^{T-1} \left(\mathbb{E} \|\nabla f(\overline{X}_t)\|^2 + \left(1 - L\gamma - \frac{6L^2C_2\gamma^4L^2}{C_3} \right) \mathbb{E} \|\overline{\nabla f}(X_t)\|^2 \right) \\ & \leq \frac{2(f(0) - f^*)}{\gamma} + \frac{LT\gamma\sigma^2}{n} + \frac{12L^2C_2\gamma^2\sigma^2T}{C_3} + \frac{6L^2C_2\gamma^4L^2\sigma^2T}{nC_3} \\ & \quad + \frac{6L^2C_1\gamma^2\sigma^2}{C_3} + \frac{6L^2C_1\gamma^2\zeta_0^2}{C_3} + \frac{6L^2C_1\gamma^2}{C_3} \|\nabla f(\mathbf{0})\|^2. \end{aligned}$$

Then we have

$$\left(1 - \frac{6L^2C_1\gamma^2}{C_3} \right) \|\nabla f(\mathbf{0})\|^2 + \sum_{t=1}^{T-1} \left(\mathbb{E} \|\nabla f(\overline{X}_t)\|^2 + \left(1 - L\gamma - \frac{6L^2C_2\gamma^4L^2}{C_3} \right) \mathbb{E} \|\overline{\nabla f}(X_t)\|^2 \right)$$

$$\leq \frac{2(f(0) - f^*)}{\gamma} + \frac{LT\gamma}{n}\sigma^2 + \frac{12L^2C_2\gamma^2\sigma^2T}{C_3} + \frac{6L^2C_2\gamma^4L^2\sigma^2T}{nC_3} + \frac{6L^2C_1\gamma^2\sigma^2}{C_3} + \frac{6L^2C_1\gamma^2\zeta_0^2}{C_3}.$$

Denote

$$A_1 = 1 - \frac{6L^2C_1\gamma^2}{C_3}$$

$$A_2 = 1 - L\gamma - \frac{6L^2C_2\gamma^4L^2}{C_3},$$

it becomes

$$A_1\|\nabla f(\mathbf{0})\|^2 + \sum_{t=1}^{T-1} \left(\mathbb{E}\|\nabla f(\bar{X}_t)\|^2 + A_2\mathbb{E}\|\bar{\nabla}f(X_t)\|^2 \right)$$

$$\leq \frac{2(f(0) - f^*)}{\gamma} + \frac{LT\gamma}{n}\sigma^2 + \frac{12L^2C_2\gamma^2n\sigma^2T}{nC_3} + \frac{6L^2C_2\gamma^4L^2\sigma^2T}{nC_3} + \frac{6L^2C_1\gamma^2\sigma^2}{C_3} + \frac{6L^2C_1\gamma^2\zeta_0^2}{C_3}.$$

It completes the proof. \square

Proof to Corollary 3

Proof. From the value of γ , we obtain

$$C_2\gamma^2L^2 \leq \frac{1}{64}, \quad C_1\gamma^2L^2 \leq \frac{1}{36}.$$

Therefore

$$C_3 = 1 - 24C_2\gamma^2L^2 \geq \frac{1}{2},$$

$$A_1 = 1 - \frac{6L^2C_1\gamma^2}{C_3} \geq \frac{1}{2},$$

$$A_2 = 1 - L\gamma - \frac{6L^2C_2\gamma^4L^2}{C_3} > 0,$$

$$\gamma^2 \leq \frac{n}{nL^2 + \sigma^2T},$$

$$\gamma^4 \leq \frac{n^2}{n^2L^4 + \sigma^4T^2}.$$

Then we can remove the $\|\bar{\nabla}f(X_t)\|^2$ and $\|\nabla f(\mathbf{0})\|^2$ on the left hand side of (5) in Theorem 2, and (5) becomes

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\|\nabla f(\bar{X}_t)\|^2 \leq \frac{4(f(0) - f^*)L(8\sqrt{C_2} + 6\sqrt{C_1})}{T} + \frac{4(f(0) - f^*)\sigma}{\sqrt{Tn}}$$

$$+ \frac{2L\sigma}{\sqrt{Tn}} + \frac{48nL^2C_2\sigma^2}{nL^2 + \sigma^2T} + \frac{24L^4n\sigma^2C_2}{n^2L^4 + \sigma^4T^2}$$

$$+ \frac{24nL^2C_1\sigma^2}{T(nL^2 + \sigma^2T)} + \frac{24L^2C_1\zeta_0^2}{T(nL^2 + \sigma^2T)},$$

which completes the proof. \square