

Decentralized Online Learning: Exchanging Local Models to Track Dynamics

Anonymous Authors¹

Abstract

In this paper, we consider online learning in the decentralized setting, which is motivated by the application scenario where users want to take benefits from the data from other users, but do not want to share their private data to a third party or other users. Instead, they can only share their private prediction model, e.g., recommendation model. We study the decentralized online gradient method in which each user maintains a private model and share its private model with its neighbors (or users he/she trusts) periodically. In addition, to consider more practical scenario we allow users' interest changing over time (it means that the optimal model changes over time), unlike most online works which assume that the optimal prediction model is constant. We show that decentralized online gradient (DOG) can efficiently and effectively propagate the values in all private data without sharing them to track the dynamics of users' interest, by proving a tight dynamic regret $\mathcal{O}\left(n\sqrt{TM} + \sqrt{nTM}\sigma\right)$ for DOG where n is the number of users, T is the number of time steps, M measures the dynamics (this is, how much the users' interest changes over time), and σ measures the randomness of the private data. Empirical studies are also conducted to validate our analysis. This study indicates the possibility of a new framework of data service: all users can take benefit from their private data without sharing them.

1. Introduction

Online learning has been studied for decades of years in machine learning literatures (????????). The goal of online learning generally is to incrementally learn predictions models to minimize the sum of all the online loss functions

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

(cumulative loss), which is usually determined by a sequence of examples that arrives sequentially. To quantify the efficacy of an online learning algorithm, the community introduced a performance measure called static regret, which is the difference between the cumulative losses suffered by the online algorithm and that suffered by the best model which can observe all the loss functions. The best static regret of a sequential online convex optimization method is $\mathcal{O}(\sqrt{T})$ and $\mathcal{O}(\log T)$ for convex and strongly convex loss functions, respectively (???)

Different from traditional online learning, online learning in decentralized networks (or Decentralized Online Learning) assumes that a network of computational nodes can communicate between neighbors to solve an online learning problem, in which each computational node will receive a stream of online losses. Suppose we have n workers, among which the i -th one will receive the t -th loss $f_{i,t}$ at the t -th iteration. Then, the goal of Decentralized Online Learning usually is to minimize its static regret, which is defined as the difference between the cumulative loss over all the nodes and steps and that of the best model which knows all the loss function beforehand; Decentralized Online Learning enjoys many advantages for real-world large-scale applications. Firstly, it avoids collecting all the loss functions to one node, which will result in heavy communication cost for the network and extremely high computational cost for one node. Secondly, it can help many data providers collaborate to better minimize their cumulative loss, while at the same time protect the data privacy as much as possible.

The static regret assumes that the best model keeps unchanged during the entire learning process, however this does not hold in some real applications. For example, one's favorite style of music may change over time as his/her situation. To solve this issue, the dynamic regret is introduced, which generally measures the difference between the cumulative loss suffered by the decentralized online learning algorithm and that suffered by a dynamic sequence of models. This dynamic sequence of models can not only observe all the loss functions beforehand, but also changes over time with the amount of changes less than a budget. In this paper, we mainly prove that decentralized online gradient can achieve a dynamic regret of $\mathcal{O}\left(n\sqrt{TM} + \sqrt{nTM}\sigma\right)$ where n is the number of users, T is the number of time steps, M mea-

sure the dynamics budget, and σ measures the randomness of the private data.

Notations and definitions In the paper, we make the following notations.

- For any $i \in [n]$ and $t \in [T]$, the random variable $\xi_{i,t}$ is subject to a distribution $D_{i,t}$, that is, $\xi_{i,t} \sim D_{i,t}$. Besides, a set of random variables $\Xi_{n,T}$ and the corresponding set of distributions are defined by $\Xi_{n,T} = \{\xi_{i,t}\}_{1 \leq i \leq n, 1 \leq t \leq T}$, and $\mathcal{D}_{n,T} = \{D_{i,t}\}_{1 \leq i \leq n, 1 \leq t \leq T}$, respectively. For math brevity, we use the notation $\Xi_{n,T} \sim \mathcal{D}_{n,T}$ to represent that $\xi_{i,t} \sim D_{i,t}$ holds for any $i \in [n]$ and $t \in [T]$. \mathbb{E} represents mathematical expectation.
- For a decentralized network, we use $\mathbf{W} \in \mathbb{R}^{n \times n}$ to represent its confusion matrix. It is a symmetric doubly stochastic matrix, which implies that every element of \mathbf{W} is non-negative, $\mathbf{W}\mathbf{1} = \mathbf{1}$, and $\mathbf{1}^T \mathbf{W} = \mathbf{1}^T$. We use $\{\lambda_i\}_{i=1}^n$ with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ to represent its eigenvalues.
- ∇ represents gradient operator. $\|\cdot\|$ represents the ℓ_2 norm in default.
- \lesssim represents “less than equal up to a constant factor”.
- \mathcal{A} represents the set of all online algorithms.

2. Related work

Online learning has been studied for decades of years. The static regret of a sequential online convex optimization method can achieve $\mathcal{O}(\sqrt{T})$ and $\mathcal{O}(\log T)$ bounds for convex and strongly convex loss functions, respectively (???). Recently, both the decentralized online learning and the dynamic regret have drawn much attention due to their wide existence in the practical big data scenarios.

2.1. Decentralized online learning

Online learning in a decentralized network has been studied in (????????????). ? studies decentralized online mirror descent, and provides $\mathcal{O}(n\sqrt{nTM})$ dynamic regret. Here, n , T , and M represent the number of nodes in the network, the number of iterations, and the budget of dynamics (defined in (??)), respectively. When the Bregman divergence in the decentralized online mirror descent is chosen appropriately, the decentralized online mirror descent becomes identical to the decentralized online gradient descent. Using the same definition of dynamic regret (defined in (??)), our method obtains $\mathcal{O}(n\sqrt{TM})$ dynamic regret

for a decentralized online gradient descent, which is better than $\mathcal{O}(n\sqrt{nTM})$ in ?. The improvement of our bound benefits from a better bound of network error (see Lemma ??). ? studies decentralized online prediction, and presents $\mathcal{O}(\sqrt{nT})$ static regret. It assumes that all data, used to yielded the loss, is generated from an unknown distribution. The strong assumption is not practical in the dynamic environment, and thus limits its novelty for a general online learning task. Additionally, many decentralized online optimization methods are proposed, for example, decentralized online multi-task learning (?), decentralized online ADMM (?), decentralized online gradient descent (?), decentralized continuous-time online saddle-point method (?), decentralized online Nesterov’s primal-dual method (??), and online distributed dual averaging(?). Those previous methods are proved to yield $\mathcal{O}(\sqrt{T})$ static regret, which do not have theoretical guarantee of regret in the dynamic environment. Besides, ? provides necessary and sufficient conditions to preserve privacy for decentralized online learning methods, which is interesting to extend our method to be privacy-preserving in the future work.

2.2. Dynamic regret

Dynamic regret has been widely studied for decades of years (????????????). ? first defines the dynamic regret by (??), and then proposes an online gradient descent method. The method yields $\mathcal{O}(\sqrt{TM})$ by choosing an appropriate learning rate. The following researches achieve the sublinear dynamic regret, but extend the analysis of regret by using different reference points. For example, ?? choose the reference points $\{\mathbf{x}_t^*\}_{t=1}^T$ satisfying $\sum_{t=1}^{T-1} \|\mathbf{x}_{t+1}^* - \Phi(\mathbf{x}_t^*)\| \leq M$, where $\Phi(\mathbf{x}_t^*)$ is the predictive optimal model. When the function Φ predicts accurately, a small M is enough to bound the dynamics. The dynamic regret is thus effectively decreased. ?????? chooses the reference points $\{\mathbf{y}_t^*\}_{t=1}^T$ with $\mathbf{y}_t^* = \arg\min_{\mathbf{z} \in \mathcal{X}} f_t(\mathbf{z})$, where f_t is the loss function at the t -th iteration. ? provides a new analysis framework, which achieves $\mathcal{O}(\sqrt{TM})$ dynamic regret¹ for any given reference points. Besides, ? presents that the lower bound of the dynamic regret defined by ?? is $\Omega(\sqrt{TM})$. The previous definition of the regret, i.e., (??), is a special case of our new definition. When setting $\gamma = 1$, we achieve the state-of-the-art regret, that is, $\mathcal{O}(\sqrt{TM})$.

In some literatures, the regret in a dynamic environment is measured by the number of changes of a reference point over time. It is usually denoted by shifting regret or tracking

¹? uses the notation of “shifting regret” instead of “dynamic regret”. In the paper, we keep using “dynamic regret” as used in most previous literatures.

regret (????????). Both the shifting regret and the tracking regret can be considered as a variation of the dynamic regret, and is usually studied in the setting of “learning with expert advice”. But, the dynamic regret is usually studied in a general setting of online learning.

3. Problem formulation

Suppose that there are n users. Each user maintains a local predictive model, and only talk to his/her neighbors. Let $\mathbf{x}_{i,t}$ denote the local model for user i at iteration t . In iteration t user i applies the local model $\mathbf{x}_{i,t}$ to a function $f_{i,t}(\cdot; \xi_{i,t})$ and receives the loss $f_{i,t}(\cdot; \xi_{i,t})$. $\xi_{i,t}$ is an i.i.d. random variable in terms of i and t , charactering the *random* component in the function $f_{i,t}(\cdot; \xi_{i,t})$, while the subscripts i and t of f (as well ξ) indicate the *adversarial* component, for example, the user’s profile, location, local time, and etc. The random component in the function is usually **To Peilin: please provide some examples here.**

Communication network. Users do not want to share the information to others and can only share their private models to their neighbors (or friends). The graph is denoted by $\mathcal{G} = (\text{nodes}: [n], \text{edges}: E)$. **Chen: please use this notation to.**

Suppose we have n workers, among which the i -th one will receive the t -th loss $f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t})$ at the t -th iteration, where $\mathbf{x}_{i,t}$ is the model on the node i at the t -th iteration, and $\xi_{i,t}$ represents the randomness of the private data, which is drawn from some distribution $D_{i,t}$ for the i -th node at the time t . If we do not consider the dynamic setup, all nodes want to make their models converge to the best model \mathbf{x}^* , which is defined by

$$\mathbf{x}^* := \arg\min_{\mathbf{x}} \sum_{i=1}^n \sum_{t=1}^T \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} f_{i,t}(\mathbf{x}; \xi_{i,t}).$$

However, the static optimal model \mathbf{x}^* does not exist in the dynamic environment. For example, when we want to conduct music recommendation to a user, his/her preference to music may change over time as his/her situation. Thus, the optimal model should be time-varying over time. It leads to the dynamics of the optimal recommendation model. Therefore, we allow the optimal model can change as the times goes, and denote them by $\{\mathbf{x}_t^*\}_{t=1}^T$. Formally, given the budget of dynamics M , the feasible reference points $\{\mathbf{x}_t^*\}_{t=1}^T$ satisfy

$$\sum_{t=1}^{T-1} \|\mathbf{x}_{t+1}^* - \mathbf{x}_t^*\| \leq M,$$

where M is used to control how quickly the optimal model can change. When $M = 0$, all \mathbf{x}_t^* s are same, and it degenerates to the static online learning problem. When the

dynamic environment changes significantly, M becomes large to model the dynamics. Let us take an example to explain the dynamics.

Denote $\mathcal{L}_M^T = \left\{ \{\mathbf{z}_t\}_{t=1}^T : \sum_{t=1}^{T-1} \|\mathbf{z}_{t+1} - \mathbf{z}_t\| \leq M \right\}$. We define a new dynamic regret as follows.

Definition 1. For any a decentralized online algorithm $A \in \mathcal{A}$, we define its dynamic regret \mathcal{R}_T^A by

$$\mathcal{R}_T^A := \mathbb{E}_{\mathbf{z}_t \sim \mathcal{L}_M^T} \sum_{i=1}^n \sum_{t=1}^T (f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - f_{i,t}(\mathbf{x}_t^*; \xi_{i,t})), \quad (1)$$

where $\{\mathbf{x}_t^*\}_{t=1}^T$ is defined by

$$\{\mathbf{x}_t^*\}_{t=1}^T := \arg\min_{\{\mathbf{z}_t\}_{t=1}^T \in \mathcal{L}_M^T} \sum_{i=1}^n \sum_{t=1}^T \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} f_{i,t}(\mathbf{z}_t; \xi_{i,t}).$$

Here, $\mathbf{x}_{i,t}$ is the model played by an online algorithm A at the t -th iteration. $\xi_{i,t}$ represents the stochastic factor in the local model, which is drawn from the distribution $D_{i,t}$ for the i -th node at the time t . It is the major difference between our new regret and the previous regret. When we do not exploit the relation among local models, e.g. $\mathbf{x}_{i,t}$, our new dynamic regret (??) de-generates to the classic definition (??). Specifically, for any a decentralized online algorithm $A \in \mathcal{A}$, the previous dynamic regret $\tilde{\mathcal{R}}_T^A$ is usually defined by

$$\tilde{\mathcal{R}}_T^A := \sum_{i=1}^n \sum_{t=1}^T (f_{i,t}(\mathbf{x}_{i,t}) - f_{i,t}(\mathbf{x}_t^*)), \quad (2)$$

subject to $\sum_{t=1}^{T-1} \|\mathbf{x}_{t+1}^* - \mathbf{x}_t^*\| \leq M$. In (??), the classic online learning in a decentralized network treats all the data as the adversary data in default, which ignores the potential relation among local models $\{\mathbf{x}_{i,t}\}_{i=1}^n$ at time t . But, those local models are not independent in many practical scenarios. For example, users’ preference to music may be impacted by a popular trend in the Internet at the same time. In the decentralized online learning, every node shares its private model to neighbours, and the regret caused by the stochastic part of data would be decreased effectively. Exploiting the potential relation among local models is helpful to reduce the regret, and learn a better model, which is varified by the theoretical results in Section ??.

4. Decentralized online gradient method

In the section, we first present the decentralized online gradient method, and then prove that it leads to $\mathcal{O}(n\sqrt{TM} + \sqrt{nTM}\sigma)$ dynamic regret.

Algorithm 1 DOG: Decentralized Online Gradient method.

Require: The learning rate η , number of iterations T , and the confusion matrix \mathbf{W} . $\mathbf{x}_{i,1} = \mathbf{0}$ for any $i \in [n]$.

- 1: **for** $t = 1, 2, \dots, T$ **do**
- 2: **For the** i -**th node with** $i \in [n]$:
- 3: Predict $\mathbf{x}_{i,t}$.
- 4: Observe the loss function $f_{i,t}$, and suffer loss $f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t})$.
- 5: **Update:**
- 6: Query a gradient $\nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t})$.
- 7: $\mathbf{x}_{i,t+1} = \sum_{j=1}^n \mathbf{W}_{i,j} \mathbf{x}_{j,t} - \eta \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t})$.
- 8: **end for**

4.1. Algorithm

The Decentralized Online Gradient method, namely DOG, is presented in Algorithm ???. This algorithm works iteration by iteration. At each iteration, every node needs to collect local models, e.g., $\mathbf{x}_{i,t}$, from its neighbours, and compute a weighted sum as $\sum_{j=1}^n \mathbf{W}_{i,j} \mathbf{x}_{j,t}$. Then, the weighted sum is updated by an online gradient descent step. In addition, we denote $\bar{\mathbf{x}}_t = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{i,t}$ to facilitate the theoretical analysis. We can verify that $\bar{\mathbf{x}}_{t+1} = \bar{\mathbf{x}}_t - \frac{\eta}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t})$ (see Lemma ??).

4.2. Theoretical analysis

Denote

$$F_{i,t}(\cdot) := \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} f_{i,t}(\cdot; \xi_{i,t}).$$

Assumption 1. We make following assumptions to analyze the dynamic regret theoretically.

- For any $i \in [n]$, $t \in [T]$, and \mathbf{x} , there exist constants G and σ^2 such that

$$\mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \|\nabla f_{i,t}(\mathbf{x}; \xi_{i,t})\|^2 \leq G,$$

and

$$\mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \|\nabla f_{i,t}(\mathbf{x}; \xi_{i,t}) - \nabla F_{i,t}(\mathbf{x})\|^2 \leq \sigma^2.$$

- For given vectors \mathbf{x} and \mathbf{y} , we assume $\|\mathbf{x} - \mathbf{y}\|^2 \leq R$.
- For any $i \in [n]$ and $t \in [T]$, we assume the function $f_{i,t}$ is convex, and has L -Lipschitz gradient.
- Given a symmetric doubly stochastic matrix \mathbf{W} , and a constant ρ with $\rho := \max\{|\lambda_2(\mathbf{W})|, |\lambda_n(\mathbf{W})|\}$, we assume $\rho < 1$.

The bound of dynamic regret yielded by Algorithm ??? is presented in the following theorem.

Theorem 1. Denote constants C_0 , and C_1 by

$$C_0 := \frac{L + 2\eta L^2 + 4L^2\eta}{(1 - \rho)^2} + 2L.$$

Using Assumption ??, and choosing $\eta > 0$ in Algorithm ??, we have

$$\begin{aligned} & \mathbb{E}_{n,T \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \sum_{i=1}^n f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - f_{i,t}(\mathbf{x}_t^*; \xi_{i,t}) \\ & \leq 20\eta T n G + \eta T \sigma^2 + C_0 n T \eta^2 G + \frac{n}{2\eta} (4\sqrt{R}M + R). \end{aligned}$$

By choosing an approximate learning rate η , we obtain sublinear regret as follows.

Corollary 1. Using Assumption ??, and choosing

$$\eta = \sqrt{\frac{(1 - \rho) (nM\sqrt{R} + nR)}{nTG + T\sigma^2}}$$

in Algorithm ??, we have

$$\begin{aligned} & \mathcal{R}_T^{\text{DOG}} \\ & \lesssim n \sqrt{T (M + \sqrt{R}) G} + \sqrt{nT (M + \sqrt{R}) \sigma^2} \\ & \quad + \frac{n (M + \sqrt{R})}{1 - \rho} + \sqrt{\frac{TM (n^2 G + n \sigma^2)}{1 - \rho}} \\ & \quad + \sqrt{\frac{T (n^2 G + n \sigma^2)}{1 - \rho}}. \end{aligned} \quad (3)$$

First, Corollary ?? shows that the dynamic regret of DOG is sublinear. Second, we would like make some comments on the effects of different parameters on the dynamic regret. The regret becomes large with the increase of the budget of dynamics M . When $n = 1$ and $\rho = 0$, the dynamic regret is $\mathcal{O}(\sqrt{TM} + \sqrt{T})$, which is tight in the case of $n = 1$ (?).

When $\rho < 1$, the regret $\mathcal{R}_T^{\text{DOG}}$ has $\sqrt{nTM\sigma^2}$ dependence on σ^2 , instead of $\sqrt{n^2TM\sigma^2}$. It benefits from the communication among nodes in the decentralized setting. Since every node shares its model with its neighbours, the variance of the average of stochastic gradients $\frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t})$ is decreased to be $\frac{\sigma^2}{n}$, thus eventually reducing the regret caused by the stochastic part of data. Additionally, the regret is affected by the topology of the network, which is measured by ρ with $0 \leq \rho < 1$. For a fully connected network², $\rho = 0$, then the regret is better than those for other topologies.

²When a network is fully connected, a decentralized method de-generates to a centralized method.

4.3. Discussions with previous work

Dependence on n . ? investigates the dynamic regret $\tilde{\mathcal{R}}_T^{\text{DOG}}$ by using DOG, and provide the following sublinear regret.

Theorem 2 (Implied by Theorem 3 and Corollary 4 in ?). *Use Assumption ??, and choose $\eta = \sqrt{\frac{(1-\rho)M}{T}}$ in Algorithm ??. The dynamic regret $\tilde{\mathcal{R}}_T^{\text{DOG}}$ is bounded by $\mathcal{O}\left(n^{\frac{3}{2}}\sqrt{\frac{MT}{1-\rho}}\right)$.*

As illustrated in Theorem ??, ? has provided a $\mathcal{O}\left(n\sqrt{nTM}\right)$ regret for DOG by using the previous dynamic regret defined in (?). Compared with the result in ?, our regret enjoys the state-of-the-art dependence on T and M , and meanwhile improves the dependence on n . This improvement is achieved by a better bound on the difference between $\mathbf{x}_{i,t}$ and $\bar{\mathbf{x}}_t$ ³.

Lemma 1. *Using Assumption ??, and setting $\eta > 0$ in Algorithm ??, we have*

$$\mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{i=1}^n \sum_{t=1}^T \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 \leq \frac{nT\eta^2 G}{(1-\rho)^2}.$$

Dependence on σ^2 . Previous researches (???) view all data as the adversary data, ignoring the potential relations among local models. They usually assume gradient of the loss function $\nabla f_{i,t}$ is bounded, e.g., $\|\nabla f_{i,t}(\mathbf{x}; \xi_{i,t}, \xi_{i,t})\|^2 \leq G$, which implies $\|\nabla h_t(\mathbf{x}; \xi_{i,t})\|^2 \leq G$, and $\mathbb{E}_{\xi_{i,t} \sim \mathcal{D}_{i,t}} \|\nabla h_t(\mathbf{x}; \xi_{i,t})\|^2 \leq \sigma^2 + G$ according to Lemma ??.

Lemma 2. *Assume $\|\nabla h_t(\mathbf{x}; \xi_{i,t})\|^2 \leq G$. It implies*

$$\mathbb{E}_{\xi_{i,t} \sim \mathcal{D}_{i,t}} \|\nabla h_t(\mathbf{x}; \xi_{i,t})\|^2 \leq \sigma^2 + G.$$

Using this assumption in previous analysis frameworks, the regret $\mathcal{R}_T^{\text{DOG}}$ has the same dependence on both G and σ^2 even in the static environment. However, our new analysis shows that the regret $\mathcal{R}_T^{\text{DOG}}$ has $\sqrt{n\sigma^2}$ dependence on σ^2 , and $\sqrt{n^2 G}$ dependence on G . The reason is that the variance of the average of stochastic gradients, i.e., $\nabla h_t(\cdot, \xi_{i,t})$ with $i \in [n]$, is decreased effectively when every node shares its local model to others.

5. Empirical studies

For simplicity, in the experiments we only consider online logistic regression with squared ℓ_2 norm regularization, i.e., $f_{i,t}(\mathbf{x}; \xi_{i,t}) = \log(1 + \exp(-\mathbf{y}_{i,t} \mathbf{A}_{i,t}^T \mathbf{x})) + \frac{\gamma}{2} \|\mathbf{x}\|^2$, where $\gamma = 10^{-3}$ is a given hyper-parameter. Under this setting, we compare the proposed Decentralized Online Gradient

³? denotes $\|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|$ by “network error”.

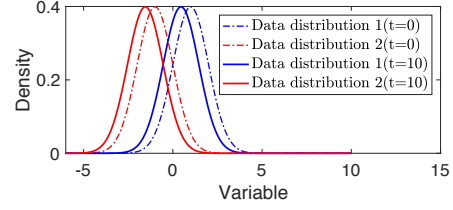


Figure 1. An illustration of the dynamics caused by the time-varying distributions of data. Data distributions 1 and 2 satisfy $N(1 + \sin(t), 1)$ and $N(-1 + \sin(t), 1)$, respectively. Suppose we want to conduct classification between data drawn from distributions 1 and 2, respectively. The optimal classification model should change over time.

method (DOG) and the Centralized Online Gradient method (COG).

M is fixed as 10 to determine the space of reference points. The learning rate η is tuned to be optimal for each dataset separately. We evaluate the learning performance by measuring the average loss $\frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t})$, instead of the dynamic regret $\mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{i=1}^n \sum_{t=1}^T (f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - f_{i,t}(\mathbf{x}_t^*))$, since the optimal reference point $\{\mathbf{x}_t^*\}_{t=1}^T$ is the same for both DOG and COG.

5.1. Datasets

To test the proposed algorithm, we utilized a toy dataset and three real-world datasets, whose details are presented as follows.

Synthetic Data For the i -th node, a data matrix $\mathbf{A}_i \in \mathbb{R}^{10 \times T}$ is generated, s.t. $\mathbf{A}_i = 0.1\hat{\mathbf{A}}_i + 0.9\tilde{\mathbf{A}}_i$, where $\tilde{\mathbf{A}}_i$ represents the adversary part of data, and $\hat{\mathbf{A}}_i$ represents the stochastic part of data. Specifically, elements of $\hat{\mathbf{A}}_i$ is uniformly sampled from the interval $[-0.5 + \sin(i), 0.5 + \sin(i)]$. Note that $\hat{\mathbf{A}}_i$ and $\hat{\mathbf{A}}_j$ with $i \neq j$ are drawn from different distributions. $\hat{\mathbf{A}}_{i,t}$ is generated according to $\mathbf{y}_{i,t} \in \{1, -1\}$ which is generated uniformly. When $\mathbf{y}_{i,t} = 1$, $\hat{\mathbf{A}}_{i,t}$ is generated by sampling from a time-varying distribution $N((1 + 0.5 \sin(t)) \cdot \mathbf{1}, \mathbf{I})$. When $\mathbf{y}_{i,t} = -1$, $\hat{\mathbf{A}}_{i,t}$ is generated by sampling from another time-varying distribution $N((-1 + 0.5 \sin(t)) \cdot \mathbf{1}, \mathbf{I})$. Due to this correlation, $\mathbf{y}_{i,t}$ can be considered as the label of the instance $\hat{\mathbf{A}}_{i,t}$. The above dynamics of time-varying distributions are illustrated in Figure ??, which shows the change of the optimal learning model over time and the importance of studying the dynamic regret.

Real Data Three real public datasets are *room-occupancy*⁴,

⁴<https://archive.ics.uci.edu/ml/datasets/Occupancy+Detection+>

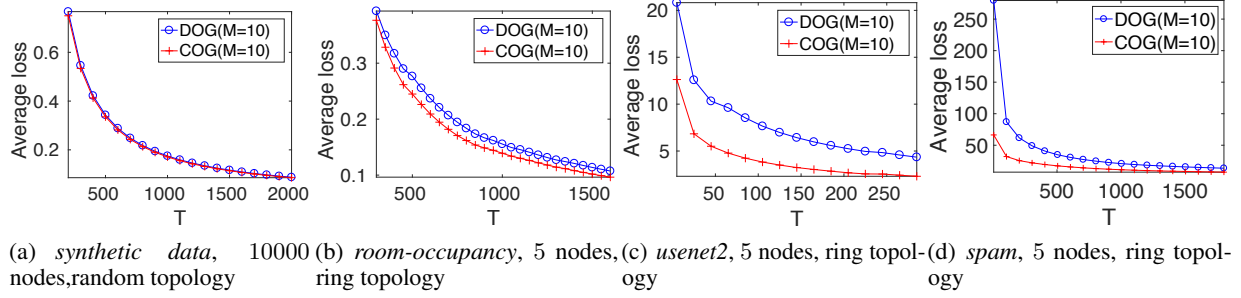


Figure 2. The average loss yielded by DOG is comparable to that yielded by COG.

*usenet2*⁵, and *spam*⁶. *room-occupancy* is a time-series dataset, which is from a natural dynamic environment. Both *usenet2* and *spam* are “concept drift” (?) datasets, for which the optimal model changes over time. Before conducting experiments, we conduct clustering for all instances, and then place all the instances within a cluster in a node to guarantee the distribution of instances for every node is different.

5.2. Results

First, Figure ?? summarizes the performance of DOG compared with COG on all the datasets. For the synthetic dataset, we simulated a decentralized network consisting of 10000 nodes, where every node is randomly connected with other 15 nodes. For the three real datasets, we simulated a network consisting of 5 nodes. In these networks, the nodes are connected by a ring topology. Under these settings, we can observe that both DOG and COG are effective for the online learning tasks on all the datasets, while DOG achieves slightly worse performance.

Second, Figure ?? summarizes the effect of the network size on the performance of DOG. We change the number of nodes from 5000 to 10000 on the synthetic dataset, and from 5 to 20 on the real datasets. The synthetic dataset is tested by using the random topology, and those real datasets are tested by using the ring topology. Figure ?? draws the curves of average loss over time steps. We observe that the average loss curves are mostly overlapped with different nodes. It shows that DOG is robust to the network size (or number of users), which validates our theory, that is, the average regret does not increase with the number of nodes. Furthermore, we observe that the average loss becomes large with the increase of the variance of stochastic data, which validates our theoretical result nicely.

Third, Figure ?? shows the effect of the topology of the network on the performance of DOG, where five differ-

ρ	NC	FC	Ring	WS(1)	Ws(0.5)
synthetic data	1	0	0.99	0.37	0.58
real data	1	0	0.96	0.83	0.76

Table 1. ρ in different topologies used in our experiment. “NC” represents the *No connected* topology, “FC” represents the *Fully connected* topology, and “WS” represents the *WattsStrogatz* topology.

ent topologies are used. Besides the ring topology, the *No connected* topology means there are no edges in the network, and every node does not share its local model to others. The *Fully connected* topology means all nodes are connected, where DOG de-generates to be COG. The topology *WattsStrogatz* represents a Watts-Strogatz small-world graph, for which we can use a parameter to control the number of random edges (set as 0.5 and 1 in this paper). The result shows *Fully connected* enjoys the best performance, because that $\rho = 0$ for it while $\rho > 0$ for other topologies. Specifically, ρ in those topologies is presented in Table ?. A small ρ leads to a good performance of DOG, which validates our theoretical result nicely.

6. Conclusion

We investigate a new online learning problem in a decentralized network, where the loss incurs by both adversary and stochastic data. We provide a new analysis framework, which achieves sublinear regret. Extensive empirical studies verify the theoretical result.

References

- D. Adamskiy, W. M. Koolen, A. Chernov, and V. Vovk. A closer look at adaptive regret. *Journal of Machine Learning Research*, 17(23):1–21, 2016.
- M. Akbari, B. Ghahesifard, and T. Linder. Distributed online convex optimization on time-varying directed graphs. *IEEE Transactions on Control of Network Systems*, 4(3): 417–428, Sep. 2017.

⁵http://mlkd.csd.auth.gr/concept_drift.html

⁶http://mlkd.csd.auth.gr/concept_drift.html

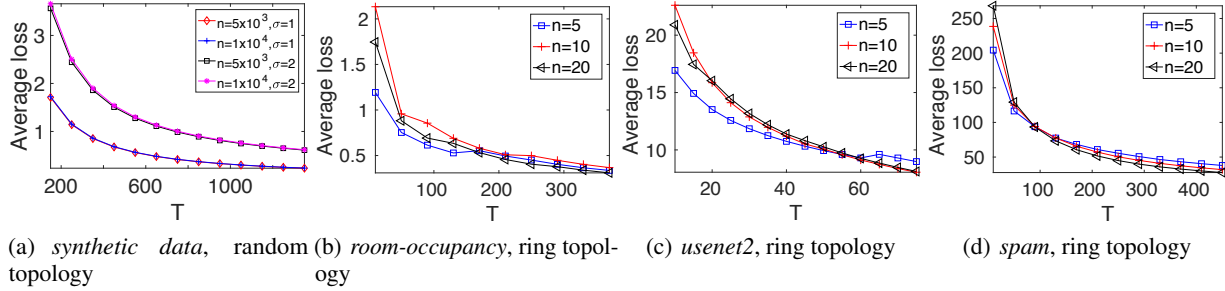


Figure 3. The average loss yielded by DOG is insensitive to the network size.

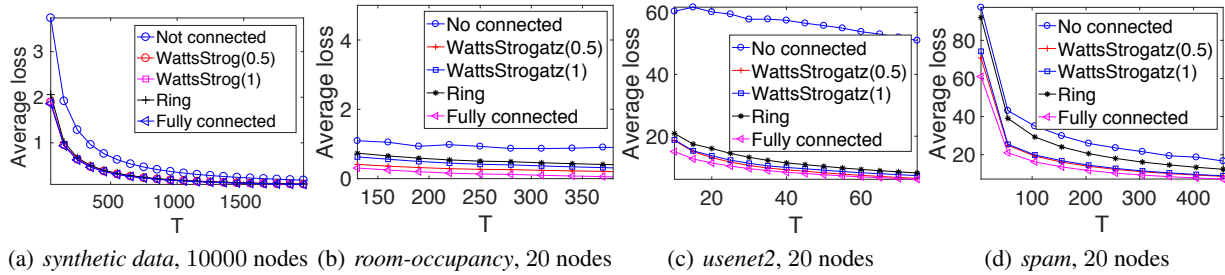


Figure 4. The average loss yielded by DOG is insensitive to the topology of the network.

- F. Bach and V. Perchet. Highly-Smooth Zero-th Order Online Optimization Vianney Perchet. *arXiv.org*, May 2016.
- A. S. Bedi, P. Sarma, and K. Rajawat. Tracking moving agents via inexact online gradient descent algorithm. *IEEE Journal of Selected Topics in Signal Processing*, 12(1):202–217, Feb 2018.
- A. A. Benczúr, L. Kocsis, and R. Pálóvics. Online Machine Learning in Big Data Streams. *CoRR*, 2018.
- S. Bubeck. Introduction to online optimization, December 2011.
- N. Cesa-Bianchi, P. Gaillard, G. Lugosi, and G. Stoltz. Mirror Descent Meets Fixed Share (and feels no regret). In *NIPS 2012*, page Paper 471, 2012.
- N. Chen, G. Goel, and A. Wierman. Smoothed Online Convex Optimization in High Dimensions via Online Balanced Descent. *arXiv.org*, Mar. 2018.
- A. György and C. Szepesvári. Shifting regret, mirror descent, and matrices. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48, ICML’16*, pages 2943–2951. JMLR.org, 2016.
- A. György, T. Linder, and G. Lugosi. Tracking the Best of Many Experts. *Proceedings of Conference on Learning Theory (COLT)*, 2005.
- A. György, T. Linder, and G. Lugosi. Efficient tracking of large classes of experts. *IEEE Transactions on Information Theory*, 58(11):6709–6725, Nov 2012.
- E. C. Hall and R. Willett. Dynamical Models and tracking regret in online convex programming. In *Proceedings of International Conference on International Conference on Machine Learning (ICML)*, 2013.
- E. C. Hall and R. M. Willett. Online Convex Optimization in Dynamic Environments. *IEEE Journal of Selected Topics in Signal Processing*, 9(4):647–662, 2015.
- E. Hazan. Introduction to online convex optimization. *Foundations and Trends in Optimization*, 2(3-4):157–325, 2016.
- M. Herbster and M. K. Warmuth. Tracking the best expert. *Machine Learning*, 32(2):151–178, Aug 1998.
- N. Ho-Nguyen and F. Kilinc-Karzan. Exploiting Problem Structure in Optimization under Uncertainty via Online Convex Optimization. *arXiv.org*, (3):741–35, Sept. 2017.
- S. Hosseini, A. Chapman, and M. Mesbahi. Online distributed optimization via dual averaging. In *52nd IEEE Conference on Decision and Control*, pages 1484–1489, Dec 2013.
- A. Jadbabaie, A. Rakhlin, S. Shahrampour, and K. Sridharan. Online Optimization : Competing with Dynamic

- Comparators. In *Proceedings of International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 398–406, 2015.
- K.-S. Jun, F. Orabona, S. Wright, and R. Willett. Improved strongly adaptive online learning using coin betting. In A. Singh and J. Zhu, editors, *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 54, pages 943–951, 20–22 Apr 2017.
- M. Kamp, M. Boley, D. Keren, A. Schuster, and I. Sharfman. Communication-efficient distributed online prediction by dynamic model synchronization. In *Proceedings of the 2014th European Conference on Machine Learning and Knowledge Discovery in Databases - Volume Part I, ECMLPKDD’14*, pages 623–639, Berlin, Heidelberg, 2014. Springer-Verlag.
- I. Katakis, G. Tsoumakas, and I. Vlahavas. Tracking recurring contexts using ensemble classifiers: An application to email filtering. *Knowledge and Information Systems*, 22(3):371–391, 2010.
- A. Koppel, S. Paternain, C. Richard, and A. Ribeiro. Decentralized online learning with kernels. *IEEE Transactions on Signal Processing*, 66(12):3240–3255, June 2018.
- S. Lee, A. Ribeiro, and M. M. Zavlanos. Distributed continuous-time online optimization using saddle-point methods. In *2016 IEEE 55th Conference on Decision and Control (CDC)*, pages 4314–4319, Dec 2016.
- S. Lee, A. Nedić, and M. Raginsky. Coordinate dual averaging for decentralized online optimization with nonseparable global objectives. *IEEE Transactions on Control of Network Systems*, 5(1):34–44, March 2018.
- M. Mohri and S. Yang. Competing with automata-based expert sequences. In A. Storkey and F. Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84, pages 1732–1740, 09–11 Apr 2018.
- A. Mokhtari, S. Shahrampour, A. Jadbabaie, and A. Ribeiro. Online optimization in dynamic environments: Improved regret rates for strongly convex problems. In *Proceedings of IEEE Conference on Decision and Control (CDC)*, pages 7195–7201. IEEE, 2016.
- J. Mourtada and O.-A. Maillard. Efficient tracking of a growing number of experts. *arXiv.org*, Aug. 2017.
- A. Nedić, S. Lee, and M. Raginsky. Decentralized online optimization with global objectives and local communication. In *2015 American Control Conference (ACC)*, pages 4497–4503, July 2015.
- M. J. Neely and H. Yu. Online Convex Optimization with Time-Varying Constraints. *arXiv.org*, Feb. 2017.
- F. Orabona, J. Luo, and B. Caputo. Multi Kernel Learning with Online-Batch Optimization. *Journal of Machine Learning Research*, 2012.
- S. Paternain and A. Ribeiro. Online Learning of Feasible Strategies in Unknown Environments. *arXiv.org*, Apr. 2016.
- S. Shahrampour and A. Jadbabaie. Distributed online optimization in dynamic environments using mirror descent. *IEEE Transactions on Automatic Control*, 63(3):714–725, March 2018.
- S. Shalev-Shwartz. Online Learning and Online Convex Optimization. *Foundations and Trends® in Machine Learning*, 4(2):107–194, 2012.
- H. Tang, S. Gan, C. Zhang, T. Zhang, and J. Liu. Communication Compression for Decentralized Training. *arXiv.org*, Mar. 2018.
- C.-Y. Wei, Y.-T. Hong, and C.-J. Lu. Tracking the best expert in non-stationary stochastic environments. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Proceedings of Advances in Neural Information Processing Systems*, pages 3972–3980, 2016.
- H.-F. Xu, Q. Ling, and A. Ribeiro. Online learning over a decentralized network through admm. *Journal of the Operations Research Society of China*, 3(4):537–562, Dec 2015.
- F. Yan, S. Sundaram, S. V. N. Vishwanathan, and Y. Qi. Distributed autonomous online learning: Regrets and intrinsic privacy-preserving properties. *IEEE Transactions on Knowledge and Data Engineering*, 25(11):2483–2493, Nov 2013.
- T. Yang, L. Zhang, R. Jin, and J. Yi. Tracking Slowly Moving Clairvoyant - Optimal Dynamic Regret of Online Learning with True and Noisy Gradient. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, 2016.
- C. Zhang, P. Zhao, S. Hao, Y. C. Soh, B. S. Lee, C. Miao, and S. C. H. Hoi. Distributed multi-task classification: a decentralized online learning approach. *Machine Learning*, 107(4):727–747, Apr 2018a.
- L. Zhang, T. Yang, J. Yi, R. Jin, and Z.-H. Zhou. Improved Dynamic Regret for Non-degenerate Functions. In *Proceedings of Neural Information Processing Systems (NIPS)*, 2017a.

L. Zhang, T. Yang, rong jin, and Z.-H. Zhou. Dynamic regret of strongly adaptive methods. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, pages 5882–5891, 10–15 Jul 2018b.

W. Zhang, P. Zhao, W. Zhu, S. C. H. Hoi, and T. Zhang. Projection-free distributed online learning in networks. In D. Precup and Y. W. Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, pages 4054–4062, International Convention Centre, Sydney, Australia, 06–11 Aug 2017b.

Y. Zhao, S. Qiu, and J. Liu. Proximal Online Gradient is Optimum for Dynamic Regret. *CoRR*, cs.LG, 2018.

M. Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of International Conference on Machine Learning (ICML)*, pages 928–935, 2003.

Appendix

Proof to Theorem ??:

Proof. From the regret definition, we have

$$\begin{aligned}
& \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \frac{1}{n} \sum_{i=1}^n f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - f_{i,t}(\mathbf{x}_t^*; \xi_{i,t}) \\
& \leq \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \frac{1}{n} \sum_{i=1}^n \langle \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \mathbf{x}_{i,t} - \mathbf{x}_t^* \rangle \\
& = \underbrace{\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \frac{1}{n} \sum_{i=1}^n (\langle \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle + \langle \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1} \rangle)}_{I_1(t)} \\
& \quad + \underbrace{\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\langle \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \bar{\mathbf{x}}_{t+1} - \mathbf{x}_t^* \right\rangle}_{I_2(t)}.
\end{aligned}$$

Now, we begin to bound $I_1(t)$.

$$I_1(t) = \left(\underbrace{\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \frac{1}{n} \sum_{i=1}^n \langle \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle}_{J_1(t)} + \underbrace{\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\langle \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1} \right\rangle}_{J_2(t)} \right).$$

For $J_1(t)$, we have

$$\begin{aligned}
& J_1(t) \\
& = \frac{1}{n} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \sum_{i=1}^n \langle \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle \\
& = \frac{1}{n} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \sum_{i=1}^n \langle \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - \nabla F_{i,t}(\bar{\mathbf{x}}_t), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle + \frac{1}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \sum_{i=1}^n \langle \nabla F_{i,t}(\bar{\mathbf{x}}_t), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle \\
& = \frac{1}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \sum_{i=1}^n \langle \nabla F_{i,t}(\mathbf{x}_{i,t}) - \nabla F_{i,t}(\bar{\mathbf{x}}_t), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle + \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \left\langle \nabla F_{i,t}(\bar{\mathbf{x}}_t), \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \right\rangle \\
& \stackrel{\textcircled{1}}{\leq} \frac{L}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2.
\end{aligned}$$

① holds due to $F_{i,t}$ has L -Lipschitz gradients, and $\bar{\mathbf{x}}_t = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{i,t}$.

For $J_2(t)$, we have

$$\begin{aligned}
& J_2(t) \\
& = \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\langle \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1} \right\rangle \\
& \leq \frac{\eta}{2} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2 + \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2
\end{aligned}$$

$$\begin{aligned}
 &\leq \frac{\eta}{2} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\| \frac{1}{n} \sum_{i=1}^n (\nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - \nabla F_{i,t}(\mathbf{x}_{i,t}) + \nabla F_{i,t}(\mathbf{x}_{i,t})) \right\|^2 + \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 \\
 &\leq \eta \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\| \frac{1}{n} \sum_{i=1}^n (\nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - \nabla F_{i,t}(\mathbf{x}_{i,t})) \right\|^2 + \eta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \left\| \frac{1}{n} \sum_{i=1}^n \nabla F_{i,t}(\mathbf{x}_{i,t}) \right\|^2 \\
 &\quad + \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 \\
 &\stackrel{\textcircled{1}}{\leq} \frac{\eta}{n} \sigma^2 + \eta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \left\| \frac{1}{n} \sum_{i=1}^n (\nabla F_{i,t}(\mathbf{x}_{i,t}) - \nabla F_{i,t}(\bar{\mathbf{x}}_t) + \nabla F_{i,t}(\bar{\mathbf{x}}_t)) \right\|^2 + \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 \\
 &\leq \frac{\eta}{n} \sigma^2 + 2\eta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \left\| \frac{1}{n} \sum_{i=1}^n (\nabla F_{i,t}(\mathbf{x}_{i,t}) - \nabla F_{i,t}(\bar{\mathbf{x}}_t)) \right\|^2 \\
 &\quad + 2\eta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \|\nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2 + \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 \\
 &\leq \frac{\eta}{n} \sigma^2 + \frac{2\eta}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \sum_{i=1}^n \|\nabla F_{i,t}(\mathbf{x}_{i,t}) - \nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2 \\
 &\quad + 2\eta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \|\nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2 + \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 \\
 &\stackrel{\textcircled{2}}{\leq} \frac{\eta}{n} \sigma^2 + \frac{2\eta L^2}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + 2\eta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \|\nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2 + \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2.
 \end{aligned}$$

① holds due to

$$\begin{aligned}
 &\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\| \frac{1}{n} \sum_{i=1}^n (\nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - \nabla F_{i,t}(\mathbf{x}_{i,t})) \right\|^2 \\
 &= \frac{1}{n^2} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \left(\sum_{i=1}^n \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \|\nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - \nabla F_{i,t}(\mathbf{x}_{i,t})\|^2 \right) \\
 &\quad + \frac{1}{n^2} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \left(2 \sum_{i=1}^n \sum_{j=1, j \neq i}^n \left\langle \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - \nabla F_{i,t}(\mathbf{x}_{i,t}), \mathbb{E}_{\xi_{j,t} \sim D_{j,t}} \nabla f_{j,t}(\mathbf{x}_{j,t}; \xi_{j,t}) - \nabla F_{j,t}(\mathbf{x}_{j,t}) \right\rangle \right) \\
 &= \frac{1}{n^2} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \sum_{i=1}^n \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \|\nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - \nabla F_{i,t}(\mathbf{x}_{i,t})\|^2 + 0 \\
 &\leq \frac{1}{n} \sigma^2.
 \end{aligned}$$

② holds due to $F_{i,t}$ has L Lipschitz gradients.

Therefore, we obtain

$$\begin{aligned}
 &I_1(t) \\
 &= (J_1(t) + J_2(t)) \\
 &= \left(\frac{L}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + \frac{\eta}{n} \sigma^2 + \frac{2\eta L^2}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 \right) \\
 &\quad + \left(2\eta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \|\nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2 + \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 \right) \\
 &\leq \left(\frac{L}{n} + \frac{2\eta L^2}{n} \right) \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + 2\eta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \|\nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2
 \end{aligned}$$

$$+ \frac{\eta\sigma^2}{n} + \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2.$$

Therefore, we have

$$\begin{aligned} \sum_{t=1}^T I_1(t) &\leq \left(\frac{L}{n} + \frac{2\eta L^2}{n} \right) \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \sum_{i=1}^n \sum_{t=1}^T \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + 2\eta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \sum_{t=1}^T \|\nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2 \\ &\quad + \frac{T\eta\sigma^2}{n} + \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \sum_{t=1}^T \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2. \end{aligned}$$

Now, we begin to bound $I_2(t)$. Recall that the update rule is

$$\mathbf{x}_{i,t+1} = \sum_{j=1}^n \mathbf{W}_{ij} \mathbf{x}_{j,t} - \eta \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}).$$

According to Lemma ??, we have

$$\bar{\mathbf{x}}_{t+1} = \bar{\mathbf{x}}_t - \eta \left(\frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right). \quad (4)$$

Denote a new auxiliary function $\phi(\mathbf{z})$ as

$$\phi(\mathbf{z}) = \left\langle \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \mathbf{z} \right\rangle + \frac{1}{2\eta} \|\mathbf{z} - \bar{\mathbf{x}}_t\|^2.$$

It is trivial to verify that (??) satisfies the first-order optimality condition of the optimization problem: $\min_{\mathbf{z} \in \mathbb{R}^d} \phi(\mathbf{z})$, that is,

$$\nabla \phi(\bar{\mathbf{x}}_{t+1}) = \mathbf{0}.$$

We thus have

$$\begin{aligned} \bar{\mathbf{x}}_{t+1} &= \operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^d} \phi(\mathbf{z}) \\ &= \operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^d} \left\langle \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \mathbf{z} \right\rangle + \frac{1}{2\eta} \|\mathbf{z} - \bar{\mathbf{x}}_t\|^2. \end{aligned}$$

Furthermore, denote a new auxiliary variable $\bar{\mathbf{x}}_\tau$ as

$$\bar{\mathbf{x}}_\tau = \bar{\mathbf{x}}_{t+1} + \tau (\mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1}),$$

where $0 < \tau \leq 1$. According to the optimality of $\bar{\mathbf{x}}_{t+1}$, we have

$$\begin{aligned} 0 &\leq \phi(\bar{\mathbf{x}}_\tau) - \phi(\bar{\mathbf{x}}_{t+1}) \\ &= \left\langle \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \bar{\mathbf{x}}_\tau - \bar{\mathbf{x}}_{t+1} \right\rangle + \frac{1}{2\eta} \left(\|\bar{\mathbf{x}}_\tau - \bar{\mathbf{x}}_t\|^2 - \|\bar{\mathbf{x}}_{t+1} - \bar{\mathbf{x}}_t\|^2 \right) \\ &= \left\langle \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \tau (\mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1}) \right\rangle + \frac{1}{2\eta} \left(\|\bar{\mathbf{x}}_{t+1} + \tau (\mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1}) - \bar{\mathbf{x}}_t\|^2 - \|\bar{\mathbf{x}}_{t+1} - \bar{\mathbf{x}}_t\|^2 \right) \\ &= \left\langle \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \tau (\mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1}) \right\rangle + \frac{1}{2\eta} \left(\|\tau (\mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1})\|^2 + 2 \langle \tau (\mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1}), \bar{\mathbf{x}}_{t+1} - \bar{\mathbf{x}}_t \rangle \right). \end{aligned}$$

Note that the above inequality holds for any $0 < \tau \leq 1$. Divide τ on both sides, and we have

$$\begin{aligned}
 I_2(t) &= \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\langle \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \bar{\mathbf{x}}_{t+1} - \mathbf{x}_t^* \right\rangle \\
 &\leq \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left(\lim_{\tau \rightarrow 0^+} \tau \|(\mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1})\|^2 + 2 \langle \mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1}, \bar{\mathbf{x}}_{t+1} - \bar{\mathbf{x}}_t \rangle \right) \\
 &= \frac{1}{\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \langle \mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1}, \bar{\mathbf{x}}_{t+1} - \bar{\mathbf{x}}_t \rangle \\
 &= \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left(\|\mathbf{x}_t^* - \bar{\mathbf{x}}_t\|^2 - \|\mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1}\|^2 - \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 \right). \tag{5}
 \end{aligned}$$

Besides, we have

$$\begin{aligned}
 &\|\mathbf{x}_{t+1}^* - \bar{\mathbf{x}}_{t+1}\|^2 - \|\mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1}\|^2 \\
 &= \|\mathbf{x}_{t+1}^*\|^2 - \|\mathbf{x}_t^*\|^2 - 2 \langle \bar{\mathbf{x}}_{t+1}, -\mathbf{x}_t^* + \mathbf{x}_{t+1}^* \rangle \\
 &= (\|\mathbf{x}_{t+1}^*\| - \|\mathbf{x}_t^*\|) (\|\mathbf{x}_{t+1}^*\| + \|\mathbf{x}_t^*\|) - 2 \langle \bar{\mathbf{x}}_{t+1}, -\mathbf{x}_t^* + \mathbf{x}_{t+1}^* \rangle \\
 &\leq \|\mathbf{x}_{t+1}^* - \mathbf{x}_t^*\| (\|\mathbf{x}_{t+1}^*\| + \|\mathbf{x}_t^*\|) + 2 \|\bar{\mathbf{x}}_{t+1}\| \|\mathbf{x}_{t+1}^* - \mathbf{x}_t^*\| \\
 &\leq 4\sqrt{R} \|\mathbf{x}_{t+1}^* - \mathbf{x}_t^*\|.
 \end{aligned}$$

The last inequality holds due to our assumption, that is, $\|\mathbf{x}_{t+1}^*\| = \|\mathbf{x}_{t+1}^* - \mathbf{0}\| \leq \sqrt{R}$, $\|\mathbf{x}_t^*\| = \|\mathbf{x}_t^* - \mathbf{0}\| \leq \sqrt{R}$, and $\|\bar{\mathbf{x}}_{t+1}\| = \|\bar{\mathbf{x}}_{t+1} - \mathbf{0}\| \leq \sqrt{R}$.

Thus, telescoping $I_2(t)$ over $t \in [T]$, we have

$$\begin{aligned}
 \sum_{t=1}^T I_2(t) &\leq \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \left(4\sqrt{R} \sum_{t=1}^T \|\mathbf{x}_{t+1}^* - \mathbf{x}_t^*\| + \|\bar{\mathbf{x}}_1^* - \bar{\mathbf{x}}_1\|^2 - \|\bar{\mathbf{x}}_T^* - \bar{\mathbf{x}}_{T+1}\|^2 \right) - \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 \\
 &\leq \frac{1}{2\eta} \left(4\sqrt{R}M + R \right) - \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2.
 \end{aligned}$$

Here, M the budget of the dynamics, which is defined in (??).

Combining those bounds of $I_1(t)$, and $I_2(t)$ together, we finally obtain

$$\begin{aligned}
 &\mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \sum_{i=1}^n f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - f_{i,t}(\mathbf{x}_t^*; \xi_{i,t}) \\
 &\leq n \sum_{t=1}^T (I_1(t) + I_2(t)) \\
 &\leq \left(\frac{L}{n} + \frac{2\eta L^2}{n} \right) \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{i=1}^n \sum_{t=1}^T \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + 2\eta \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \|\nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2 + \frac{T\eta\sigma^2}{n} + \frac{n}{2\eta} (4\sqrt{R}M + R) \\
 &\stackrel{\textcircled{1}}{\leq} \eta T \sigma^2 + 4n \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T (F_{i,t}(\bar{\mathbf{x}}_t) - F_{i,t}(\bar{\mathbf{x}}_{t+1})) + (L + 2\eta L^2 + 4L^2\eta) \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 \\
 &\quad + 4n \left(4T\eta G + \frac{TGL\eta^2}{2} \right) + \frac{n}{2\eta} (4\sqrt{R}M + R) \\
 &\stackrel{\textcircled{2}}{\leq} \eta T \sigma^2 + 4n \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T (F_{i,t}(\bar{\mathbf{x}}_t) - F_{i,t}(\bar{\mathbf{x}}_{t+1})) + (L + 2\eta L^2 + 4L^2\eta) \frac{nT\eta^2 G}{(1-\rho)^2} \\
 &\quad + 4n \left(4T\eta G + \frac{TGL\eta^2}{2} \right) + \frac{n}{2\eta} (4\sqrt{R}M + R)
 \end{aligned}$$

$$\stackrel{\textcircled{3}}{\leq} \eta T \sigma^2 + 4nT\eta G + (L + 2\eta L^2 + 4L^2\eta) \frac{nT\eta^2 G}{(1-\rho)^2} + 4n \left(4T\eta G + \frac{TGL\eta^2}{2} \right) + \frac{n}{2\eta} (4\sqrt{R}M + R).$$

① holds due to Lemma ?? . That is, we have

$$\begin{aligned} & \frac{\eta}{2} \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{t=1}^T \|\nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2 \\ & \leq \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T (F_{i,t}(\bar{\mathbf{x}}_t) - F_{i,t}(\bar{\mathbf{x}}_{t+1})) + 4T\eta G + \frac{L^2\eta}{n} \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{t=1}^T \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 + \frac{TGL\eta^2}{2}. \end{aligned}$$

② holds due to Lemma ??

$$\mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{i=1}^n \sum_{t=1}^T \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 \leq \frac{nT\eta^2 G}{(1-\rho)^2}.$$

③ holds due to

$$\begin{aligned} & \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} (F_{i,t}(\bar{\mathbf{x}}_t) - F_{i,t}(\bar{\mathbf{x}}_{t+1})) \\ & \leq \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \langle \nabla F_{i,t}(\bar{\mathbf{x}}_t), \bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1} \rangle \\ & = \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\langle \nabla F_{i,t}(\bar{\mathbf{x}}_t), \frac{\eta}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\rangle \\ & \leq \eta \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left(\frac{1}{2} \|\nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2 + \frac{1}{2} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2 \right) \\ & \leq \eta \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left(\frac{1}{2} \|\nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2 + \frac{1}{2n} \sum_{i=1}^n \|\nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t})\|^2 \right) \\ & \leq \eta G. \end{aligned}$$

Re-arranging items, we have

$$\begin{aligned} & \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \sum_{i=1}^n f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - f_{i,t}(\mathbf{x}_t^*; \xi_{i,t}) \\ & \leq 20\eta TnG + \eta T \sigma^2 + \left(\frac{L + 2\eta L^2 + 4L^2\eta}{(1-\rho)^2} + 2L \right) nT\eta^2 G + \frac{n}{2\eta} (4\sqrt{R}M + R). \end{aligned}$$

It completes the proof. □

Lemma 3. Using Assumption ??, and setting $\eta > 0$ in Algorithm ??, we have

$$\begin{aligned} & \frac{\eta}{2} \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{t=1}^T \|\nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2 \\ & \leq \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T (F_{i,t}(\bar{\mathbf{x}}_t) - F_{i,t}(\bar{\mathbf{x}}_{t+1})) + 4T\eta G + \frac{L^2\eta}{n} \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{t=1}^T \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 + \frac{TGL\eta^2}{2}. \end{aligned} \tag{6}$$

Proof. We have

$$\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} F_{i,t}(\bar{\mathbf{x}}_{t+1})$$

$$\begin{aligned}
 &\leq \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} F_{i,t}(\bar{\mathbf{x}}_t) + \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \langle \nabla F_{i,t}(\bar{\mathbf{x}}_t), \bar{\mathbf{x}}_{t+1} - \bar{\mathbf{x}}_t \rangle + \frac{L}{2} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_{t+1} - \bar{\mathbf{x}}_t\|^2 \\
 &= \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} F_{i,t}(\bar{\mathbf{x}}_t) + \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\langle \nabla F_{i,t}(\bar{\mathbf{x}}_t), -\frac{\eta}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\rangle + \frac{L}{2} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\| \frac{\eta}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2 \\
 &= \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} F_{i,t}(\bar{\mathbf{x}}_t) + \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \left\langle \nabla F_{i,t}(\bar{\mathbf{x}}_t), -\frac{\eta}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\rangle + \frac{L}{2} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\| \frac{\eta}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2.
 \end{aligned} \tag{7}$$

Besides, we have

$$\begin{aligned}
 &\mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \left\langle \nabla F_{i,t}(\bar{\mathbf{x}}_t), -\frac{\eta}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\rangle \\
 &= \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \left(\left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) - \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2 - \left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) \right\|^2 - \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2 \right) \\
 &\leq \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \left(\left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) - \frac{1}{n} \sum_{i=1}^n (\nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) + \nabla F_{i,t}(\mathbf{x}_{i,t})) \right\|^2 \right) - \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) \right\|^2 \\
 &\leq \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \left(2 \left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) - \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2 + 2 \left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) - \frac{1}{n} \sum_{i=1}^n \nabla F_{i,t}(\mathbf{x}_{i,t}) \right\|^2 \right) \\
 &\quad - \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) \right\|^2 \\
 &\leq \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \left(2 \left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) - \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2 + \frac{2}{n} \sum_{i=1}^n \left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) - \nabla F_{i,t}(\mathbf{x}_{i,t}) \right\|^2 \right) \\
 &\quad - \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) \right\|^2 \\
 &\leq \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \left(2 \left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) - \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2 + \frac{2L^2}{n} \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 \right) - \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) \right\|^2 \\
 &\leq \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \left(4 \left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) \right\|^2 + 4 \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2 + \frac{2L^2}{n} \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 \right) - \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) \right\|^2 \\
 &\stackrel{\textcircled{1}}{\leq} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \left(8G + \frac{2L^2}{n} \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 \right) - \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \frac{\eta}{2} \left\| \nabla F_{i,t}(\bar{\mathbf{x}}_t) \right\|^2.
 \end{aligned} \tag{8}$$

① holds due to

$$\mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2 \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \left\| \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2 \leq G.$$

Recall that

$$\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\| \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2 \leq G. \tag{9}$$

Substituting (??) and (??) into (??), and telescoping $t \in [T]$, we obtain

$$\mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T F_{i,t}(\bar{\mathbf{x}}_{t+1})$$

$$\begin{aligned}
 &\leq \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} F_{i,t}(\bar{\mathbf{x}}_t) + \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \left\langle \nabla F_{i,t}(\bar{\mathbf{x}}_t), -\frac{\eta}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\rangle + \frac{L}{2} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\| \frac{\eta}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2 \\
 &\leq \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} F_{i,t}(\bar{\mathbf{x}}_t) + \left(\mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \left(8G + \frac{2L^2}{n} \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 \right) - \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \|\nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2 \right) + \frac{GL\eta^2}{2} \\
 &= \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} F_{i,t}(\bar{\mathbf{x}}_t) + \left(4\eta G + \frac{L^2\eta}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 - \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{t-1}} \frac{\eta}{2} \|\nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2 \right) + \frac{GL\eta^2}{2}.
 \end{aligned}$$

Telescoping over $t \in [T]$, we have

$$\begin{aligned}
 &\frac{\eta}{2} \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{t=1}^T \|\nabla F_{i,t}(\bar{\mathbf{x}}_t)\|^2 \\
 &\leq \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T (F_{i,t}(\bar{\mathbf{x}}_t) - F_{i,t}(\bar{\mathbf{x}}_{t+1})) + 4T\eta G + \frac{L^2\eta}{n} \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{t=1}^T \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 + \frac{TGL\eta^2}{2}.
 \end{aligned} \tag{10}$$

It completes the proof. \square

Lemma 4. Denote $\bar{\mathbf{x}}_t = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{i,t}$. We have

$$\bar{\mathbf{x}}_{t+1} = \bar{\mathbf{x}}_t - \eta \left(\frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right).$$

Proof. Denote by

$$\begin{aligned}
 \mathbf{X}_t &= [\mathbf{x}_{1,t}, \mathbf{x}_{2,t}, \dots, \mathbf{x}_{n,t}] \in \mathbb{R}^{d \times n}, \\
 \mathbf{G}_t &= [\nabla f_{1,t}(\mathbf{x}_{1,t}; \zeta_{1,t}, \xi_{1,t}), \nabla f_{2,t}(\mathbf{x}_{2,t}; \zeta_{2,t}, \xi_{2,t}), \dots, \nabla f_{n,t}(\mathbf{x}_{n,t}; \zeta_{n,t}, \xi_{n,t})] \in \mathbb{R}^{d \times n}.
 \end{aligned}$$

Recall that

$$\mathbf{x}_{i,t+1} = \sum_{j=1}^n \mathbf{W}_{ij} \mathbf{x}_{j,t} - \eta \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}).$$

Equivalently, we re-formulate the update rule as

$$\mathbf{X}_{t+1} = \mathbf{X}_t \mathbf{W} - \eta \mathbf{G}_t.$$

Since the confusion matrix \mathbf{W} is doubly stochastic, we have

$$\mathbf{W} \mathbf{1} = \mathbf{1}.$$

Thus, we have

$$\begin{aligned}
 \bar{\mathbf{x}}_{t+1} &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{i,t+1} \\
 &= \mathbf{X}_{t+1} \frac{\mathbf{1}}{n} \\
 &= \mathbf{X}_t \mathbf{W} \frac{\mathbf{1}}{n} - \eta \mathbf{G}_t \frac{\mathbf{1}}{n} \\
 &= \mathbf{X}_t \frac{\mathbf{1}}{n} - \eta \mathbf{G}_t \frac{\mathbf{1}}{n} \\
 &= \bar{\mathbf{x}}_t - \eta \left(\frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right).
 \end{aligned}$$

It completes the proof. \square

Lemma 5 (Lemma 5 in (?)). For any matrix $\mathbf{X}_t \in \mathbb{R}^{d \times n}$, decompose the confusion matrix \mathbf{W} as $\mathbf{W} = \sum_{i=1}^n \lambda_i \mathbf{v}_i \mathbf{v}_i^T = \mathbf{P} \mathbf{\Lambda} \mathbf{P}^T$, where $\mathbf{P} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n] \in \mathbb{R}^{n \times n}$, \mathbf{v}_i is the normalized eigenvector of λ_i . $\mathbf{\Lambda}$ is a diagonal matrix, and λ_i be its i -th element. We have

$$\|\mathbf{X}_t \mathbf{W}^t - \mathbf{X}_t \mathbf{v}_1 \mathbf{v}_1^T\|_F^2 \leq \|\rho^t \mathbf{X}_t\|_F^2,$$

where $\rho = \max\{|\lambda_2(\mathbf{W})|, |\lambda_n(\mathbf{W})|\}$.

Lemma 6 (Lemma 6 in (?)). Given two non-negative sequences $\{a_t\}_{t=1}^\infty$ and $\{b_t\}_{t=1}^\infty$ that satisfying

$$a_t = \sum_{s=1}^t \rho^{t-s} b_s,$$

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

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

Proof to Lemma ??:

Proof. Recall that

$$\mathbf{x}_{i,t+1} = \sum_{j=1}^n \mathbf{W}_{ij} \mathbf{x}_{j,t} - \eta \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}),$$

and according to Lemma ??, we have

$$\bar{\mathbf{x}}_{t+1} = \bar{\mathbf{x}}_t - \eta \left(\frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right).$$

Denote by

$$\mathbf{X}_t = [\mathbf{x}_{1,t}, \mathbf{x}_{2,t}, \dots, \mathbf{x}_{n,t}] \in \mathbb{R}^{d \times n},$$

$$\mathbf{G}_t = [\nabla f_{1,t}(\mathbf{x}_{1,t}; \zeta_{1,t}, \xi_{1,t}), \nabla f_{2,t}(\mathbf{x}_{2,t}; \zeta_{2,t}, \xi_{2,t}), \dots, \nabla f_{n,t}(\mathbf{x}_{n,t}; \zeta_{n,t}, \xi_{n,t})] \in \mathbb{R}^{d \times n}.$$

By letting $\mathbf{x}_{i,1} = \mathbf{0}$ for any $i \in [n]$, the update rule is re-formulated as

$$\mathbf{X}_{t+1} = \mathbf{X}_t \mathbf{W} - \eta \mathbf{G}_t = - \sum_{s=1}^t \eta \mathbf{G}_s \mathbf{W}^{t-s}.$$

Similarly, denote $\bar{\mathbf{G}}_t = \frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t})$, and we have

$$\bar{\mathbf{x}}_{t+1} = \bar{\mathbf{x}}_t - \eta \left(\frac{1}{n} \sum_{i=1}^n \nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right) = - \sum_{s=1}^t \eta \bar{\mathbf{G}}_s.$$

Therefore, we obtain

$$\begin{aligned} & \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 \\ & \stackrel{\textcircled{1}}{=} \sum_{i=1}^n \left\| \sum_{s=1}^{t-1} \eta \bar{\mathbf{G}}_s - \eta \mathbf{G}_s \mathbf{W}^{t-s-1} \mathbf{e}_i \right\|^2 \end{aligned}$$

$$\begin{aligned}
 & \stackrel{\textcircled{2}}{=} \left\| \sum_{s=1}^{t-1} \eta \mathbf{G}_s \mathbf{v}_1 \mathbf{v}_1^\top - \eta \mathbf{G}_s \mathbf{W}^{t-s-1} \right\|_F^2 \\
 & \stackrel{\textcircled{3}}{\leq} \left(\eta \rho^{t-s-1} \left\| \sum_{s=1}^{t-1} \mathbf{G}_s \right\|_F \right)^2 \\
 & \leq \left(\sum_{s=1}^{t-1} \eta \rho^{t-s-1} \|\mathbf{G}_s\|_F \right)^2.
 \end{aligned}$$

① holds due to \mathbf{e}_i is a unit basis vector, whose i -th element is 1 and other elements are 0s. ② holds due to $\mathbf{v}_1 = \frac{1}{\sqrt{n}}$. ③ holds due to Lemma ??.

Thus, we have

$$\begin{aligned}
 & \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{i=1}^n \sum_{t=1}^T \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 \\
 & \leq \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \left(\sum_{s=1}^{t-1} \eta \rho^{t-s-1} \|\mathbf{G}_s\|_F \right)^2 \\
 & \stackrel{\textcircled{1}}{\leq} \frac{\eta^2}{(1-\rho)^2} \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \left(\sum_{t=1}^T \|\mathbf{G}_t\|_F^2 \right) \\
 & = \frac{\eta^2}{(1-\rho)^2} \left(\mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \sum_{i=1}^n \|\nabla f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t})\|^2 \right) \\
 & \stackrel{\textcircled{2}}{=} \frac{nT\eta^2 G}{(1-\rho)^2}.
 \end{aligned}$$

① holds due to Lemma ??.

□

Proof to Lemma ??:

Proof. We have

$$\begin{aligned}
 & \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \|\nabla h_t(\mathbf{x}; \xi_{i,t})\|^2 \\
 & = \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \left\| \nabla h_t(\mathbf{x}; \xi_{i,t}) - \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \nabla h_t(\mathbf{x}; \xi_{i,t}) + \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \nabla h_t(\mathbf{x}; \xi_{i,t}) \right\|^2 \\
 & = \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \left\| \nabla h_t(\mathbf{x}; \xi_{i,t}) - \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \nabla h_t(\mathbf{x}; \xi_{i,t}) \right\|^2 + \left\| \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \nabla h_t(\mathbf{x}; \xi_{i,t}) \right\|^2 \\
 & \quad + 2 \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \left\langle \nabla h_t(\mathbf{x}; \xi_{i,t}) - \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \nabla h_t(\mathbf{x}; \xi_{i,t}), \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \nabla h_t(\mathbf{x}; \xi_{i,t}) \right\rangle \\
 & = \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \left\| \nabla h_t(\mathbf{x}; \xi_{i,t}) - \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \nabla h_t(\mathbf{x}; \xi_{i,t}) \right\|^2 + \left\| \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \nabla h_t(\mathbf{x}; \xi_{i,t}) \right\|^2 \\
 & \leq \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \left\| \nabla h_t(\mathbf{x}; \xi_{i,t}) - \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \nabla h_t(\mathbf{x}; \xi_{i,t}) \right\|^2 + \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \|\nabla h_t(\mathbf{x}; \xi_{i,t})\|^2 \\
 & \leq \sigma^2 + G.
 \end{aligned}$$

It thus completes the proof.

□