

Gossip Online Learning: Exchanging Local Models to Track Dynamics

January 11, 2019

Abstract

1 Introduction

For any online algorithm $A \in \mathcal{A}$, the previous dynamic regret $\tilde{\mathcal{R}}_T^A$ is defined by

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

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 [Hazan, 2016, Shalev-Shwartz, 2012]. 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 [Shahrampour and Jadbabaie, 2018, Kamp et al., 2014, Koppel et al., 2018, Zhang et al., 2018a, 2017b, Xu et al., 2015, Akbari et al., 2017, Lee et al., 2016, Nedi et al., 2015, Lee et al., 2018, Benczúr et al., 2018, Yan et al., 2013]. Shahrampour and Jadbabaie [2018] studies decentralized online mirror descent, and provides $\mathcal{O}(n\sqrt{nTM})$ dynamic regret. 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. Comparing with the previous result, our method obtains $\mathcal{O}(\sqrt{nTM})$ dynamic regret (defined in (1)) for a decentralized online gradient descent.

Kamp et al. [2014] studies decentralized online prediction, and presents $\mathcal{O}(\sqrt{nT})$ static regret. It assumes that all data, used to yield the loss, is generated from an unknown distribution. The strong assumption limits its novelty for a general online learning task. Additionally, many decentralized online optimization and learning methods are proposed, for example, decentralized online multi-task learning [Zhang et al., 2018a], decentralized online ADMM [Xu et al., 2015], decentralized online sub-gradient descent [Akbari et al., 2017], decentralized continuous-time online saddle-point method [Lee et al., 2016], decentralized online Nesterov's primal-dual method [Nedi et al., 2015, Lee et al., 2018]. Those previous methods are proved to yield $\mathcal{O}(\sqrt{T})$ static regret, which do not have theoretical guarantee of regrets in the dynamic environment. Besides, Benczúr et al. [2018] reviews online learning methods for big data streams. Yan et al. [2013] provides necessary and sufficient conditions to preserve privacy for decentralized online learning methods.

2.2 Regret in dynamic environment

Dynamic regret has been widely studied for decades of years [Zinkevich, 2003, Hall and Willett, 2015, 2013, Jadbabaie et al., 2015, Yang et al., 2016, Bedi et al., 2018, Zhang et al., 2017a, Mokhtari et al., 2016, Zhang et al., 2018b, György and Szepesvári, 2016, Wei et al., 2016, Zhao et al., 2018]. Zinkevich [2003] first defines the reference points $\{\mathbf{x}_t^*\}_{t=1}^T$ satisfying (3), 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 it to different reference points. For example, Hall and Willett [2015, 2013] 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 decision variable. When the function Φ predicts accurately, a small M is enough to bound the dynamics. The dynamic regret is thus effectively decreased. Jadbabaie et al. [2015], Yang et al. [2016], Bedi et al. [2018], Zhang et al. [2017a], Mokhtari et al. [2016], Zhang et al. [2018b] chooses the reference points $\{\mathbf{y}_t^*\}_{t=1}^T$ with $\mathbf{y}_t^* = \operatorname{argmin}_{\mathbf{z} \in \mathcal{X}} f_t(\mathbf{z})$, where f_t is the loss function at the t -th iteration. György and Szepesvári [2016] provides a new analysis framework, which achieves $\mathcal{O}(\sqrt{TM})$ dynamic regret for any given reference points. Besides, Zhao et al. [2018] presents that the lower bound of the dynamic regret is $\mathcal{O}(\sqrt{TM})$. Those previous methods define the regret as (1), which is a special case of our definition. When setting $\beta = 1$, 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 regret. [Herbster and Warmuth, 1998, György et al., 2005, Gyorgy et al., 2012, György and Szepesvári, 2016, Mourtada and Maillard, 2017, Adamskiy et al., 2016, Wei et al., 2016, Cesa-Bianchi et al., 2012, Mohri and Yang, 2018, Jun et al., 2017]. 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 online setting.

3 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

$$\Xi_{n,T} = \{\xi_{i,t}\}_{1 \leq i \leq n, 1 \leq t \leq T}, \text{ and } \mathcal{D}_{n,T} = \{D_{i,t}\}_{1 \leq i \leq n, 1 \leq t \leq T},$$

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. ∂ and ∇ represent sub-gradient and gradient operators, respectively. $\|\cdot\|$ represents the ℓ_2 norm in default.

4 Problem formulation

4.1 Setup

For any online algorithm $A \in \mathcal{A}$, define its dynamic regret as

$$\mathcal{R}_T^A = \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \left(\sum_{i=1}^n \sum_{t=1}^T f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - f_{i,t}(\mathbf{x}_t^*; \xi_{i,t}) \right), \quad (2)$$

where n is the number of nodes in the decentralized network. The local loss function $f_{i,t}(\mathbf{x}; \xi_{i,t})$ is defined by

$$f_{i,t}(\mathbf{x}; \xi_{i,t}) := \beta g_{i,t}(\mathbf{x}) + (1 - \beta) h_t(\mathbf{x}; \xi_{i,t})$$

with $0 < \beta < 1$, and $\xi_{i,t}$ is a random variable drawn from an unknown distribution $D_{i,t}$. Note that $g_{i,t}$ is an adversary loss function, which is yielded by the learning model. $h_t(\cdot; \xi_{i,t})$ is a known loss function, which depends on the random variable $\xi_{i,t}$. The expectation of $h_t(\cdot; \xi_{i,t})$ is a global model, and does not depend on the i -th node.

$\{\mathbf{x}_t^*\}_{t=1}^T$ is the sequence of reference points, and

$$\{\mathbf{x}_t^*\}_{t=1}^T \in \left\{ \{\mathbf{z}_t\}_{t=1}^T : \sum_{t=1}^{T-1} \|\mathbf{z}_t - \mathbf{z}_{t+1}\| \leq M \right\}.$$

Here, M is the budget of the dynamics, that is,

$$\sum_{t=1}^{T-1} \|\mathbf{x}_{t+1}^* - \mathbf{x}_t^*\| \leq M. \quad (3)$$

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. Besides, we denote

$$H_t(\cdot) = \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} h_t(\cdot; \xi_{i,t}),$$

and

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

Recall that the previous definition of the dynamic regret is (1). Using (1), the classic online learning in a decentralized network only considers the loss function, i.e., $g_{i,t}$, incurred by the learning model on every node. Comparing with it, our definition of the dynamic regret, i.e., (2), still considers the loss function, i.e., H_t . It is incurred by a global model, which is used to let the decision variables, e.g., $\mathbf{x}_{i,t}$, have some good property in practical scenarios. We present some application scenarios to explain it in Section 4.2.

4.2 Application scenarios

To protect privacy, users prefer to placing their data in the local node, instead of providing it to a centralized server. Decentralized computing provides an alternative method to solve the problem. There is a user named Bob, who subscribes the online music recommendation service.

Online music recommendation with unreliable features. In the task, we want to decide whether to recommend some a music to Bob's mobile phone by using historical browser records of users on the Youtube. But, some values of features in those records are not reliable. For example, Alice does not want to let others know her real birthday and age. She submits random numbers for such information when signing up as an Youtube user. Note that those unreliable values, e.g., Alice's age and birthday, usually do not change, or have insignificant change over time. It can be modeled by an unknown distribution $D_{i,t}$ for the i -th user at time t . But, other reliable values, e.g., Alice's preference to music, may change over time due to time-varying trends of hot topics in the Internet. The dynamic nature of data implies that the optimal learning model should change over time. Thus, it is necessary to use dynamic regret to measure the quality of the learning model. In the case, $g_{i,t}(\mathbf{x}_{i,t})$ represents the loss incurred by those reliable features in the learning model, e.g., preference to music. $h_t(\mathbf{x}_{i,t}; \xi_{i,t})$ represents the loss incurred by those unreliable features in the learning model, e.g., age and birthday. A small β means significant attention on those unreliable features.

Suppose we use logistic regression to decide whether to recommend some a music to Bob. Without loss of generality, features corresponding to those unreliable values are denoted by the beginning s features. Given a user's behavior record $\mathbf{a}_{i,t}$ and its label $\mathbf{y}_{i,t} \in \{1, -1\}$. In the case, $g_{i,t}(\mathbf{x}) = \log \left(1 + \exp \left(-\mathbf{y}_{i,t} \mathbf{a}_{i,t}^T \hat{\mathbf{I}} \mathbf{x} \right) \right)$, where $\hat{\mathbf{I}}$ is yielded by letting the first s diagonal elements of an identity matrix be 0s. $\xi_{i,t} = \check{\mathbf{I}} \mathbf{a}_{i,t} \mathbf{y}_{i,t}^T$, and $h_t(\mathbf{x}; \xi_{i,t}) = \log \left(1 + \exp \left(-\xi_{i,t}^T \mathbf{x} \right) \right)$, where $\check{\mathbf{I}}$ is yielded by letting the last $(d - s)$ diagonal elements of an

identity matrix be 0s. Here, $\xi_{i,t}$ is drawn from an unknown distribution, that is, $\xi_{i,t} \sim D_{i,t}$, and $D_{i,t}$ usually changes insignificant over t , or does not change over t . In the case, $H_t(\mathbf{x})$ allows the decision variable \mathbf{x} to represent different models to treat the unreliable and reliable features.

Online music recommendation with user-specified privacy protection. In the task, we want to conduct online music recommendation with the user-specified privacy protection for Bob, because he wants to protect his data in the way he likes. We provide several choices for users to make a tradeoff between the accuracy of recommendation and the privacy protection. For example, when we use ϵ -differential privacy, these choices may include *strong privacy, weak accuracy* ($\epsilon = 0.01$), *medium privacy, medium accuracy* ($\epsilon = 0.05$), and *weak privacy, strong accuracy* ($\epsilon = 0.1$). Note that Bob's choice may change over time. For example, he may tolerate weak privacy protection to receive the newest song produced by his favorite player timely, but may want strong privacy protection when seeing a privacy-leaking news from a newspaper. In the case, $g_{i,t}(\mathbf{x}_{i,t})$ represents the loss incurred by the learning model. $h_t(\mathbf{x}_{i,t}; \xi_{i,t})$ represents the loss incurred by some a randomization encryption method, e.g., objective perturbation [Chaudhuri et al., 2011, Wang et al., 2017], to protect the privacy. Since Bob's preference to music may change over time, the optimal recommendation model should change over time. Thus, the dynamic regret is necessary to measure the quality of the model.

Similarly, suppose we want to learn a logistic regression model with the user-specified privacy protection. Given an instance $\mathbf{a}_{i,t} \in \mathbb{R}^d$ and its label $\mathbf{y}_{i,t} \in \{1, -1\}$. In the case, $g_{i,t}(\mathbf{x}) = \log(1 + \exp(-\mathbf{y}_{i,t} \mathbf{a}_{i,t}^T \mathbf{x}))$. We use the objective perturbation strategy [Chaudhuri et al., 2011, Wang et al., 2017] to protect the privacy. Specifically, we let $h_t(\mathbf{x}; \xi_{i,t}) = \mathbf{x}^T \xi_{i,t}$, where $\xi_{i,t}$ is random variable, whose density is

$$v(\mathbf{x}) = \frac{1}{\lambda} \exp(-\delta_{i,t} \|\mathbf{x}\|).$$

Here, λ is a normalizing constant, $\delta_{i,t}$ is a known function of $\epsilon_{i,t}$ for $\epsilon_{i,t}$ -differential privacy [Dwork and Roth, 2014]. For example, when $\delta_{i,t} = \epsilon_{i,t}$, $\lambda = ?$.

5 Algorithm

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

The decentralized online gradient method, namely DOG, is presented in Algorithm 1. At every iteration, every node needs to collect the decision variable, e.g., $\mathbf{x}_{i,t}$, from its neighbours, and then update its decision variable. Here, $\mathbf{W} \in \mathbb{R}^{n \times n}$ is the confusion matrix. It is a doubly stochastic matrix, which implies that every element of \mathbf{W} is non-negative, $\mathbf{W}\mathbf{1} = \mathbf{1}$, and $\mathbf{1}^T \mathbf{W} = \mathbf{1}^T$. Denote $\bar{\mathbf{x}}_t = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{i,t}$. We can verify that $\bar{\mathbf{x}}_{t+1} = \bar{\mathbf{x}}_t - \frac{\eta}{n} \sum_{i=1}^n \partial f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t})$ (see Lemma 3).

6 Theoretical analysis

Assumption 1. We make the following assumptions.

- For any $i \in [n]$, $t \in [T]$, and \mathbf{x} , there exists a constant G such that

$$\max \left\{ \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \|\nabla h_t(\mathbf{x}; \xi_{i,t})\|^2, \|\partial g_{i,t}(\mathbf{x})\|^2 \right\} \leq G,$$

and

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

- For any \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, but may be non-smooth. Furthermore, we assume the function H_t has L -Lipschitz gradients. In brief, $g_{i,t}$ may be non-convex, non-smooth. H_t is smooth, but may be non-convex. $f_{i,t}$ is convex, but may be non-smooth.

6.1 Main results

Theorem 1. Denote $\bar{\mathbf{x}}_t = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{i,t}$, and constants C_0 and C_1 by

$$C_0 := \frac{1}{\sqrt{\beta^2 + \eta}} + 4;$$

$$C_1 := \frac{\beta}{2\eta} + L + \frac{\sqrt{\beta^2 + \eta}}{2\eta} + 2\eta L^2 + C_0(1 - \beta)^2 L^2 \eta.$$

Using Assumption 1, and choosing $\eta > 0$ in Algorithm 1, we have

$$\begin{aligned} & \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \sum_{i=1}^n f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - f_{i,t}(\mathbf{x}_t^*; \xi_{i,t}) \\ & \leq \eta T (n\beta G + (1 - \beta)\sigma^2) + n(1 - \beta)C_0 \left(\mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T (H_t(\bar{\mathbf{x}}_t) - H_t(\bar{\mathbf{x}}_{t+1})) \right) \\ & \quad + (1 - \beta) \frac{nT\eta^2 G C_1}{(1 - \rho)^2} + n(1 - \beta)C_0 \left(4T\beta^2 \eta G + \frac{TGL\eta^2}{2} \right) + \frac{n}{2\eta} (4\sqrt{R}M + R). \end{aligned}$$

Corollary 1. Recall that

$$C_0 = \frac{1}{\sqrt{\beta^2 + \eta}} + 4.$$

Using Assumption 1, and choosing

$$\eta = \sqrt{\frac{nM}{T(n\beta G + (1 - \beta)\sigma^2)}}$$

in Algorithm 1, we have

$$\mathcal{R}_T^{\text{DOG}} \lesssim \sqrt{nMT(\beta nG + (1 - \beta)\sigma^2)} + n(1 - \beta)C_0 \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T (H_t(\bar{\mathbf{x}}_t) - H_t(\bar{\mathbf{x}}_{t+1})).$$

6.2 Connections with the previous results

7 Empirical studies

7.1 Experimental settings

We simulate a decentralized network consisting of 5 nodes. Those nodes are connected by using a ring topology. Besides, we conduct online logistic regression by using three time series datasets: *room-occupancy*¹, *online-retail*², *BeijingPM2.5*³, and a spam email dataset with the concept drift [Katakis et al., 2010]: *spam*⁴ in the decentralized network. The data distribution of those datasets may change over time in those practical scenarios, leading to the change of the optimal learning model during online learning. In those dynamic environment, the dynamic regret is practical and necessary.

- *room-occupancy*. It collects features of a room including temperature, humidity, light, and CO2 for every minute between 02/02/2015 and 02/10/2015. Label of an instance is whether the room is occupied. Our goal is to learn a classification model to make a decision whether the room is occupied by using those features.
- *online-retail*. It is an online retail dataset, which contains all transactions occurring between 01/12/2010 and 09/12/2011 for a UK-based and registered non-store online retail. We use three features, that is, *whether a transaction is cancelled*, *quantity*, and *unit price*. We need to train a binary classification model to make a decision whether a customer is coming from United Kingdom.
- *BeijingPM2.5*. It collects some weather features, e.g., temperature and pressure, and the PM2.5 data of US Embassy in Beijing hourly between 01/01/2010 and 12/31/2014. When the PM 2.5 index is larger than 100, the air quality is *bad*, otherwise, *good*. We want to train a binary classification model to make a decision whether the air quality is good according to features such as temperature and pressure.
- *spam*. Every instance in the dataset is an email, where the frequency of every word in the dictionary is collected. But, the distribution of words changes over time, which is denoted by *concept drift* [Katakis et al., 2010]. We want to learn a classification model to make a decision whether an email is a spam.

All values of a feature have been normalized to be zero mean and one variance. The budget of dynamics, namely M is set to be 10. For the t -th iteration, the learning rate, namely η in Algorithm 1 is set to be $\sqrt{\frac{5M}{100t}}$. As we have shown in Section 4.2, we test the dynamic regret in those three application scenarios.

7.2 Communication efficient online logistic regression

The hyper-parameter to control the communication efficiency, namely λ_t is set to be a constant 0.1.

7.3 Online logistic regression with privacy protection

7.4 Online logistic regression with unreliable features

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.

¹<https://archive.ics.uci.edu/ml/datasets/Occupancy+Detection+>

²<https://archive.ics.uci.edu/ml/datasets/Online+Retail>

³<https://archive.ics.uci.edu/ml/datasets/Beijing+PM2.5+Data>

⁴http://mlkd.csd.auth.gr/concept_drift.html

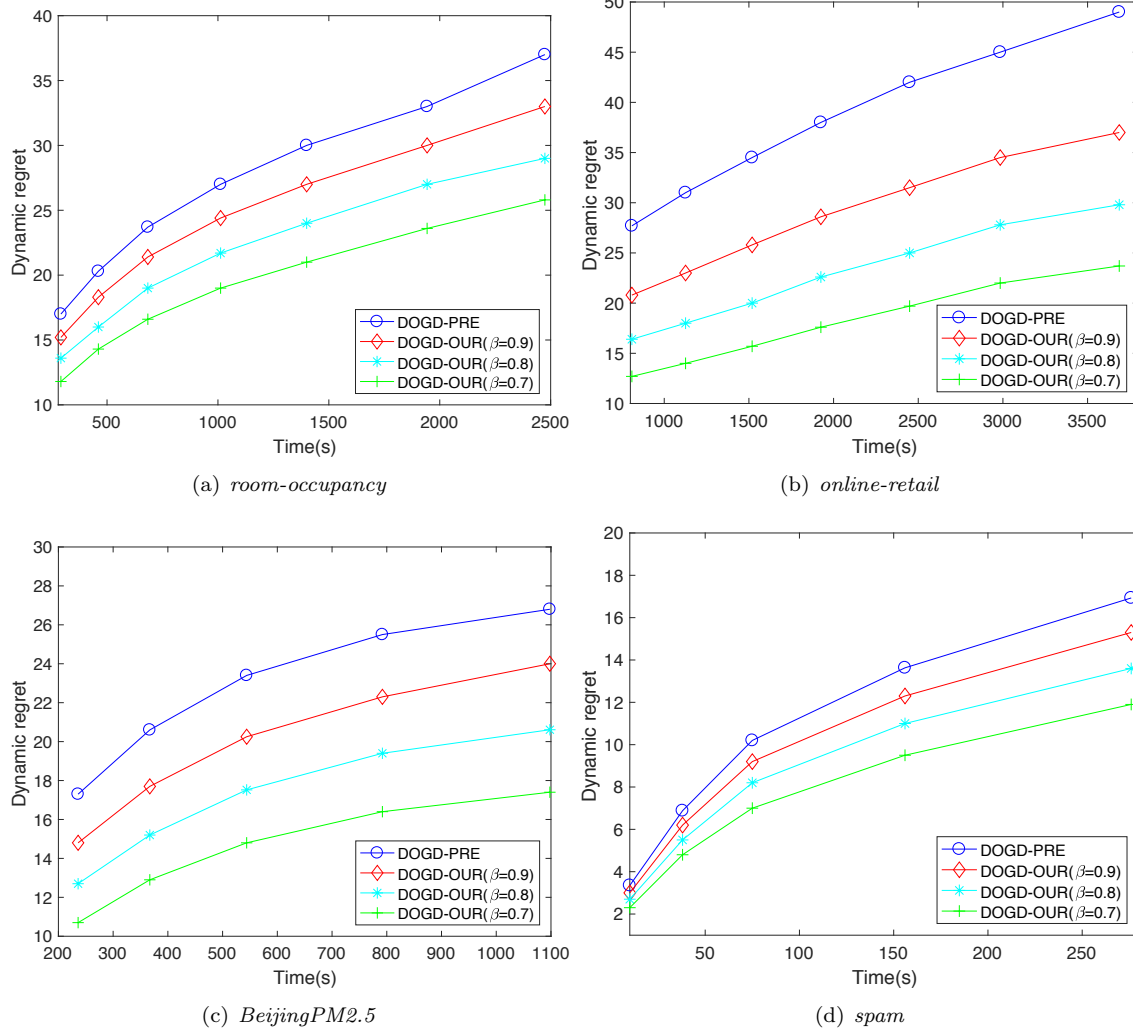


Figure 1: Comparison of the dynamic regret by using communication efficient logistic regression in the decentralized network.

- 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.
- 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.
- K. Chaudhuri, C. Monteleoni, and A. D. Sarwate. Differentially Private Empirical Risk Minimization. *Journal of Machine Learning Research*, 2011.
- C. Dwork and A. Roth. The Algorithmic Foundations of Differential Privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3-4):211–407, 2014.
- 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. Gyorgy, 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.
- 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. Nedi, 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. Nedi, 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.
- 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.
- D. Wang, M. Ye, and J. Xu. Differentially private empirical risk minimization revisited: Faster and more general. In *Advances in Neural Information Processing Systems 30*, pages 2722–2731. 2017.
- 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 1:

Proof.

$$\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 \partial f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \mathbf{x}_{i,t} - \mathbf{x}_t^* \rangle \\
& = \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \frac{1}{n} \sum_{i=1}^n \beta \langle \partial g_{i,t}(\mathbf{x}_{i,t}), \mathbf{x}_{i,t} - \mathbf{x}_t^* \rangle + (1 - \beta) \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \frac{1}{n} \sum_{i=1}^n \langle \nabla h_t(\mathbf{x}_{i,t}; \xi_{i,t}), \mathbf{x}_{i,t} - \mathbf{x}_t^* \rangle \\
& = \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \frac{1}{n} \sum_{i=1}^n \beta (\langle \partial g_{i,t}(\mathbf{x}_{i,t}), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle + \langle \partial g_{i,t}(\mathbf{x}_{i,t}), \bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1} \rangle + \langle \partial g_{i,t}(\mathbf{x}_{i,t}), \bar{\mathbf{x}}_{t+1} - \mathbf{x}_t^* \rangle) \\
& \quad + \frac{1}{n} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \sum_{i=1}^n (1 - \beta) (\langle \nabla h_t(\mathbf{x}_{i,t}; \xi_{i,t}), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle + \langle \nabla h_t(\mathbf{x}_{i,t}; \xi_{i,t}), \bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1} \rangle) \\
& \quad + \frac{1}{n} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \sum_{i=1}^n (1 - \beta) (\langle \nabla h_t(\mathbf{x}_{i,t}; \xi_{i,t}), \bar{\mathbf{x}}_{t+1} - \mathbf{x}_t^* \rangle) \\
& = \underbrace{\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \frac{1}{n} \sum_{i=1}^n \beta (\langle \partial g_{i,t}(\mathbf{x}_{i,t}), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle + \langle \partial g_{i,t}(\mathbf{x}_{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}} \frac{1}{n} \sum_{i=1}^n (1 - \beta) (\langle \nabla h_t(\mathbf{x}_{i,t}; \xi_{i,t}), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle + \langle \nabla h_t(\mathbf{x}_{i,t}; \xi_{i,t}), \bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1} \rangle)}_{I_2(t)} \\
& \quad + \underbrace{\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\langle \frac{1}{n} \sum_{i=1}^n \partial f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \bar{\mathbf{x}}_{t+1} - \mathbf{x}_t^* \right\rangle}_{I_3(t)}
\end{aligned}$$

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

$$\begin{aligned}
I_1(t) & \stackrel{\textcircled{1}}{\leq} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \frac{\beta}{n} \sum_{i=1}^n \left(\frac{\eta}{2} \|\partial g_{i,t}(\mathbf{x}_{i,t})\|^2 + \frac{1}{2\eta} \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + \frac{\eta}{2} \|\partial g_{i,t}(\mathbf{x}_{i,t})\|^2 + \frac{1}{2\eta} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 \right) \\
& \leq \beta G \eta + \frac{\beta}{2n\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + \frac{\beta}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2.
\end{aligned}$$

① holds due to $\langle \mathbf{a}, \mathbf{b} \rangle \leq \frac{\eta}{2} \|\mathbf{a}\|^2 + \frac{1}{2\eta} \|\mathbf{b}\|^2$ holds for any $\eta > 0$.

Now, we begin to bound $I_2(t)$.

$$I_2(t) = (1 - \beta) \left(\underbrace{\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \frac{1}{n} \sum_{i=1}^n \langle \nabla h_t(\mathbf{x}_{i,t}; \xi_{i,t}), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle}_{J_1(t)} + \underbrace{\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\langle \frac{1}{n} \sum_{i=1}^n \nabla h_t(\mathbf{x}_{i,t}; \xi_{i,t}), \bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1} \right\rangle}_{J_2(t)} \right).$$

For $J_1(t)$, we have

$$\begin{aligned}
J_1(t) &= \frac{1}{n} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \sum_{i=1}^n \langle \nabla h_t(\mathbf{x}_{i,t}; \xi_{i,t}), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle \\
&= \frac{1}{n} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \sum_{i=1}^n \langle \nabla h_t(\mathbf{x}_{i,t}; \xi_{i,t}) - \nabla H_t(\bar{\mathbf{x}}_t), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle + \frac{1}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \langle \nabla H_t(\bar{\mathbf{x}}_t), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle \\
&= \frac{1}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \langle \nabla H_t(\mathbf{x}_{i,t}) - \nabla H_t(\bar{\mathbf{x}}_t), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle + \frac{1}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \langle \nabla H_t(\bar{\mathbf{x}}_t), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle \\
&\stackrel{\textcircled{1}}{\leq} \frac{L}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + \frac{1}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \langle \nabla H_t(\bar{\mathbf{x}}_t), \mathbf{x}_{i,t} - \bar{\mathbf{x}}_t \rangle \\
&\stackrel{\textcircled{2}}{\leq} \frac{L}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + \frac{1}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \left(\frac{\eta}{2\nu} \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 + \frac{\nu}{2\eta} \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 \right) \\
&\leq \frac{L}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + \frac{\eta}{2\nu} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 + \frac{\nu}{2\eta n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2. \quad (4)
\end{aligned}$$

① holds due to H_t has L -Lipschitz gradients. ② holds because that $\langle \mathbf{a}, \mathbf{b} \rangle \leq \frac{\nu}{2} \|\mathbf{a}\|^2 + \frac{1}{2\nu} \|\mathbf{b}\|^2$ holds for any $\nu > 0$.

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 h_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 h_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 \\
&\leq \frac{\eta}{2} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\| \frac{1}{n} \sum_{i=1}^n (\nabla h_t(\mathbf{x}_{i,t}; \xi_{i,t}) - \nabla H_t(\mathbf{x}_{i,t}) + \nabla H_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 h_t(\mathbf{x}_{i,t}; \xi_{i,t}) - \nabla H_t(\mathbf{x}_{i,t})) \right\|^2 + \eta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \left\| \frac{1}{n} \sum_{i=1}^n \nabla H_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}_{n,t-1}} \left\| \frac{1}{n} \sum_{i=1}^n (\nabla H_t(\mathbf{x}_{i,t}) - \nabla H_t(\bar{\mathbf{x}}_t) + \nabla H_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}_{n,t-1}} \left\| \frac{1}{n} \sum_{i=1}^n (\nabla H_t(\mathbf{x}_{i,t}) - \nabla H_t(\bar{\mathbf{x}}_t)) \right\|^2 \\
&\quad + 2\eta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \|\nabla H_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}_{n,t-1}} \sum_{i=1}^n \|\nabla H_t(\mathbf{x}_{i,t}) - \nabla H_t(\bar{\mathbf{x}}_t)\|^2 \\
&\quad + 2\eta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \|\nabla H_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}$$

$$\stackrel{\textcircled{2}}{\leq} \frac{\eta}{n} \sigma^2 + \frac{2\eta L^2}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + 2\eta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 + \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2.$$

① holds due to

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

② holds due to H_t has L Lipschitz gradients.

Therefore, we obtain

$$\begin{aligned} & I_2(t) \\ &= (1 - \beta)(J_1(t) + J_2(t)) \\ &= (1 - \beta) \left(\frac{L}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + \frac{\eta}{2\nu} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 + \frac{\nu}{2\eta n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 \right) \\ & \quad + (1 - \beta) \left(\frac{\eta}{n} \sigma^2 + \frac{2\eta L^2}{n} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 \right) \\ & \quad + (1 - \beta) \left(2\eta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 + \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 \right) \\ &\leq (1 - \beta) \left(\frac{L}{n} + \frac{\nu}{2n\eta} + \frac{2\eta L^2}{n} \right) \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + \left(\frac{\eta}{2\nu} + 2\eta \right) (1 - \beta) \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 \\ & \quad + \frac{\eta(1 - \beta)\sigma^2}{n} + \frac{1 - \beta}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2. \end{aligned}$$

Combine those bounds of $I_1(t)$ and $I_2(t)$. We thus have

$$\begin{aligned} & I_1(t) + I_2(t) \\ &\leq \beta G\eta + \frac{\beta}{2n\eta} \sum_{i=1}^n \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + \frac{\beta}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 \\ & \quad + (1 - \beta) \left(\frac{L}{n} + \frac{\nu}{2n\eta} + \frac{2\eta L^2}{n} \right) \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 + \left(\frac{\eta}{2\nu} + 2\eta \right) (1 - \beta) \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 \\ & \quad + \frac{\eta(1 - \beta)\sigma^2}{n} + \frac{1 - \beta}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 \\ &= \eta \left(\beta G + \frac{(1 - \beta)\sigma^2}{n} \right) + (1 - \beta) \left(\frac{\beta}{2n\eta} + \frac{L}{n} + \frac{\nu}{2n\eta} + \frac{2\eta L^2}{n} \right) \sum_{i=1}^n \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 \end{aligned}$$

$$+ \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 + \left(\frac{\eta}{2\nu} + 2\eta \right) (1 - \beta) \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2.$$

Therefore, we have

$$\begin{aligned} & \sum_{t=1}^T (I_1(t) + I_2(t)) \\ & \leq \eta T \left(\beta G + \frac{(1 - \beta)\sigma^2}{n} \right) + (1 - \beta) \left(\frac{\beta}{2n\eta} + \frac{L}{n} + \frac{\nu}{2n\eta} + \frac{2\eta L^2}{n} \right) \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{i=1}^n \sum_{t=1}^T \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 \\ & \quad + \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 + \left(\frac{\eta}{2\nu} + 2\eta \right) (1 - \beta) \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{t=1}^T \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2. \end{aligned}$$

Now, we begin to bound $I_3(t)$. Recall that the update rule is

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

According to Lemma 3, we have

$$\bar{\mathbf{x}}_{t+1} = \bar{\mathbf{x}}_t - \eta \left(\frac{1}{n} \sum_{i=1}^n \partial f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right). \quad (5)$$

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

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

It is trivial to verify that (5) 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

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

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

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

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

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

$$\begin{aligned} I_3(t) &= \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\langle \frac{1}{n} \sum_{i=1}^n \partial f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}), \bar{\mathbf{x}}_{t+1} - \mathbf{x}_t^* \right\rangle \\ &\leq \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left(\lim_{\tau \rightarrow 0^+} \tau \|\mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1}\|^2 + 2 \langle \mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1}, \bar{\mathbf{x}}_{t+1} - \bar{\mathbf{x}}_t \rangle \right) \\ &= \frac{1}{\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \langle \mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1}, \bar{\mathbf{x}}_{t+1} - \bar{\mathbf{x}}_t \rangle \\ &= \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left(\|\mathbf{x}_t^* - \bar{\mathbf{x}}_t\|^2 - \|\mathbf{x}_t^* - \bar{\mathbf{x}}_{t+1}\|^2 - \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 \right). \end{aligned} \quad (6)$$

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 $I_3(t)$ over $t \in [T]$, we have

$$\begin{aligned} &\sum_{t=1}^T I_3(t) \\ &\leq \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \left(4\sqrt{R} \sum_{t=1}^T \|\mathbf{x}_{t+1}^* - \mathbf{x}_t^*\| + \|\bar{\mathbf{x}}_1^* - \bar{\mathbf{x}}_1\|^2 - \|\bar{\mathbf{x}}_T^* - \bar{\mathbf{x}}_{T+1}\|^2 \right) - \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2 \\ &\leq \frac{1}{2\eta} (4\sqrt{R}M + R) - \frac{1}{2\eta} \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t+1}\|^2. \end{aligned}$$

Here, M the budget of the dynamics, which is defined in (3).

Combining those bounds of $I_1(t)$, $I_2(t)$ and $I_3(t)$ together, we finally obtain

$$\begin{aligned} &\mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \sum_{i=1}^n f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) - f_t(\mathbf{x}_t^*; \xi_{i,t}) \\ &\leq n \sum_{t=1}^T (I_1(t) + I_2(t) + I_3(t)) \\ &\leq \eta T (n\beta G + (1-\beta)\sigma^2) + (1-\beta) \left(\frac{\beta}{2\eta} + L + \frac{\nu}{2\eta} + 2\eta L^2 \right) \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{i=1}^n \sum_{t=1}^T \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 \\ &\quad + n \left(\frac{\eta}{2\nu} + 2\eta \right) (1-\beta) \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{t=1}^T \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 + \frac{n}{2\eta} (4\sqrt{R}M + R) \end{aligned}$$

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

① holds due to Lemma 2. That is, we have

$$\begin{aligned}
& \frac{\eta}{2} \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{t=1}^T \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 \\
& \leq \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T (H_t(\bar{\mathbf{x}}_t) - H_t(\bar{\mathbf{x}}_{t+1})) + 4T\beta^2 \eta G + \frac{(1-\beta)^2 L^2 \eta}{n} \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{t=1}^T \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 + \frac{TGL\eta^2}{2}.
\end{aligned} \tag{7}$$

② holds due to Lemma 4

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

Letting $\nu = \sqrt{\beta^2 + \eta}$, we have

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

It completes the proof. □

Lemma 1. Using Assumption 1, we have

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

Proof.

$$\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\partial f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t})\|^2$$

$$\begin{aligned}
&= \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\beta \partial g_{i,t}(\mathbf{x}_{i,t}) + (1-\beta) \nabla h_t(\mathbf{x}_{i,t}; \xi_{i,t})\|^2 \\
&\leq \beta \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \|\partial g_{i,t}(\mathbf{x}_{i,t})\|^2 + (1-\beta) \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\nabla h_t(\mathbf{x}_{i,t}; \xi_{i,t})\|^2 \\
&\leq G.
\end{aligned}$$

It completes the proof. \square

Lemma 2. Using Assumption 1, and setting $\eta > 0$ in Algorithm 1, we have

$$\begin{aligned}
&\frac{\eta}{2} \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{t=1}^T \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 \\
&\leq \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T (H_t(\bar{\mathbf{x}}_t) - H_t(\bar{\mathbf{x}}_{t+1})) + 4T\beta^2\eta G + \frac{(1-\beta)^2 L^2 \eta}{n} \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{t=1}^T \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 + \frac{TGL\eta^2}{2}.
\end{aligned} \tag{8}$$

Proof.

$$\begin{aligned}
&\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} H_t(\bar{\mathbf{x}}_{t+1}) \\
&\leq \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} H_t(\bar{\mathbf{x}}_t) + \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \langle \nabla H_t(\bar{\mathbf{x}}_t), \bar{\mathbf{x}}_{t+1} - \bar{\mathbf{x}}_t \rangle + \frac{L}{2} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\bar{\mathbf{x}}_{t+1} - \bar{\mathbf{x}}_t\|^2 \\
&= \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} H_t(\bar{\mathbf{x}}_t) + \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\langle \nabla H_t(\bar{\mathbf{x}}_t), -\frac{\eta}{n} \sum_{i=1}^n \partial f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\rangle + \frac{L}{2} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\| \frac{\eta}{n} \sum_{i=1}^n \partial f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2 \\
&= \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} H_t(\bar{\mathbf{x}}_t) + \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \left\langle \nabla H_t(\bar{\mathbf{x}}_t), -\frac{\eta}{n} \sum_{i=1}^n \partial f_{i,t}(\mathbf{x}_{i,t}) \right\rangle + \frac{L}{2} \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \left\| \frac{\eta}{n} \sum_{i=1}^n \partial f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right\|^2.
\end{aligned} \tag{9}$$

Besides, we have

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

$$\begin{aligned}
&\leq \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \frac{\eta}{2} \left(2\beta^2 \left\| \nabla H_t(\bar{\mathbf{x}}_t) - \frac{1}{n} \sum_{i=1}^n \partial g_{i,t}(\mathbf{x}_{i,t}) \right\|^2 + \frac{2(1-\beta)^2 L^2}{n} \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 \right) - \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \frac{\eta}{2} \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 \\
&\leq \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \frac{\eta}{2} \left(4\beta^2 \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 + 4\beta^2 \left\| \frac{1}{n} \sum_{i=1}^n \partial g_{i,t}(\mathbf{x}_{i,t}) \right\|^2 + \frac{2(1-\beta)^2 L^2}{n} \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 \right) \\
&\quad - \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \frac{\eta}{2} \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 \\
&\stackrel{\textcircled{1}}{\leq} \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \frac{\eta}{2} \left(8\beta^2 G + \frac{2(1-\beta)^2 L^2}{n} \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 \right) - \mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \frac{\eta}{2} \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2. \tag{10}
\end{aligned}$$

① holds due to

$$\begin{aligned}
\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 &= \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 \\
&= \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \left\| \mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \nabla h_t(\bar{\mathbf{x}}_t; \xi_{i,t}) \right\|^2 \\
&\leq \mathbb{E}_{\Xi_{n,t-1} \sim \mathcal{D}_{n,t-1}} \left(\mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \|\nabla h_t(\bar{\mathbf{x}}_t; \xi_{i,t})\|^2 \right), \quad \forall i \in [n] \\
&\leq G,
\end{aligned}$$

and

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

According to Lemma 1, we have

$$\mathbb{E}_{\Xi_{n,t} \sim \mathcal{D}_{n,t}} \|\partial f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t})\|^2 \leq G. \tag{11}$$

Substituting (10) and (11) into (9), and telescoping $t \in [T]$, we obtain

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

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

$$\begin{aligned}
&\frac{\eta}{2} \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{t=1}^T \|\nabla H_t(\bar{\mathbf{x}}_t)\|^2 \tag{12} \\
&\leq \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T (H_t(\bar{\mathbf{x}}_t) - H_t(\bar{\mathbf{x}}_{t+1})) + 4T\beta^2 \eta G + \frac{(1-\beta)^2 L^2 \eta}{n} \mathbb{E}_{\Xi_{n,T-1} \sim \mathcal{D}_{n,T-1}} \sum_{t=1}^T \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_{i,t}\|^2 + \frac{TGL\eta^2}{2}.
\end{aligned}$$

It completes the proof. \square

Lemma 3. Denote $\bar{\mathbf{x}}_t = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{i,t}$. We have

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

Proof. Denote

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

Recall that

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

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

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

It completes the proof. □

Lemma 4. Using Assumption 1, and setting $\eta > 0$ in Algorithm 1, we have

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

Proof. Recall that

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

and according to Lemma 3, we have

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

Denote

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

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

$$\bar{\mathbf{x}}_{t+1} = \bar{\mathbf{x}}_t - \eta \left(\frac{1}{n} \sum_{i=1}^n \partial f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t}) \right) = - \sum_{s=1}^t \eta \bar{\mathbf{G}}_s. \quad (13)$$

Therefore,

$$\begin{aligned}& \sum_{i=1}^n \|\mathbf{x}_{i,t} - \bar{\mathbf{x}}_t\|^2 \\ & \stackrel{\textcircled{1}}{=} \sum_{i=1}^n \left\| \sum_{s=1}^{t-1} \eta \bar{\mathbf{G}}_s - \eta \mathbf{G}_s \mathbf{W}^{t-s-1} \mathbf{e}_i \right\|^2 \\ & \stackrel{\textcircled{2}}{=} \left\| \sum_{s=1}^{t-1} \eta \mathbf{G}_s \mathbf{v}_1 \mathbf{v}_1^T - \eta \mathbf{G}_s \mathbf{W}^{t-s-1} \right\|_F^2 \\ & \stackrel{\textcircled{3}}{\leq} \left(\eta \rho^{t-s-1} \left\| \sum_{s=1}^{t-1} \mathbf{G}_s \right\|_F \right)^2 \\ & \leq \left(\sum_{s=1}^{t-1} \eta \rho^{t-s-1} \|\mathbf{G}_s\|_F \right)^2.\end{aligned}$$

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 5.

Thus, we have

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

① holds due to Lemma 6. ② holds due to Lemma 1.

□

Lemma 5 (Appeared in Lemma 5 in [Tang et al., 2018]). 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 (Appeared in Lemma 6 in [Tang et al., 2018]). 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

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