
Online Meta-Learning on Non-convex Setting

Zhenxun Zhuang
Boston University
Boston, MA 02215
zxxzhuang@bu.edu

Kezi Yu
IQVIA
Plymouth Meeting, PA 19462
k.yu@us.imshealth.com

Songtao Lu
University of Minnesota
Minneapolis, MN 55455
lus@umn.edu

Lucas Glass
IQVIA
Plymouth Meeting, PA 19462
Lucas.Glass@iqvia.com

Yunlong Wang
IQVIA
Plymouth Meeting, PA 19462
Yunlong.Wang@iqvia.com

Abstract

The ability of continual learning over new tasks is critical for human, and is also highly desirable for modern artificial intelligent systems. The online meta-learning framework is designed for this scenario. It combines two distinct research areas: meta-learning which tries to extract prior knowledge from existing tasks for fast learning of future tasks, and online-learning which focuses on the sequential setting in which problems are revealed one by one. In this paper, we extend the original framework from convex to non-convex setting, and introduce the local regret as the alternative performance measure. We then apply this framework to stochastic settings, and show it enjoys a logarithmic local regret, is robust to any hyperparameter initialization, and empirically outperforms traditional methods.

1 Introduction

In recent years, high-capacity machine learning models, such as deep neural networks [LBH15], have achieved remarkable successes in various domains [SHM⁺16, RDGF16, AAA⁺16]. However, domains where data is scarce remain a big challenge as those models' ability to learn and generalize rely heavily on the amount of training data. On the contrary, humans can learn new skills and concepts very efficiently from just a few experiences. For example, Lake et al. [LUTG17] observed that on playing the Atari game Frostbite, to reach the performance achieved by a human player after 2 hours' training, a modern learning model needs over 300 hours' experience. The key reason here is that when encountering a new task, learning algorithms start completely from scratch; while humans are typically armed with plenty of prior knowledge accumulated from past experience which may share overlapping structures with the current task, and thus can enable efficient learning of the new task.

Meta-learning [NM92, TP12, VBL⁺16] was designed to mimic this human ability. A meta-learning algorithm is first given a set of meta-training tasks assumed to be drawn from some distribution, and attempts to extract prior knowledge applicable to all tasks in the form of a meta-learner. This meta-learner is then evaluated on an unseen task, usually assumed to be drawn from a similar distribution as the one used for training. Although meta-learning has developed rapidly in recent years, it typically assumes all meta-training tasks are available together as a batch, which doesn't capture the sequential setting of continual lifelong learning in which new tasks are revealed one after another.

Meanwhile, online Learning [CBL06] specifically tackles the sequential setting. At each round t , the algorithm makes a choice x_t , and then suffers a loss $f_t(x_t)$ revealed by a potentially adversarial environment. The standard objective is to minimize the *regret*, the difference between the cumulative

Algorithm 1 Online Meta-Learning

- 1: **Input:** \mathbf{w}_1 , a loss function $\ell(\cdot)$, a local adapter strategy $U(\cdot)$, an online learning algorithm \mathcal{A}
 - 2: **for** $t = 1, 2, \dots$ **do**
 - 3: Encounter a new task: \mathcal{T}_t
 - 4: Receive training data for current task: \mathcal{D}_t^{tr}
 - 5: Adapt the meta-learner to current task: $\hat{\mathbf{w}}_t = U(\mathbf{w}_t, \mathcal{D}_t^{tr})$
 - 6: Receive test data for current task: \mathcal{D}_t^{ts}
 - 7: Suffer a loss: $\ell_t(\mathbf{w}_t) \triangleq \ell(\hat{\mathbf{w}}_t, \mathcal{D}_t^{ts}) = \mathbb{E}_{\mathbf{x}, y \sim \mathcal{D}_t^{ts}} [\ell(\hat{\mathbf{w}}_t, \mathbf{x}; y)]$
 - 8: Update the meta-learner: $\mathbf{w}_{t+1} = \mathcal{A}(\mathbf{w}_1, \ell_1(\mathbf{w}_1), \dots, \ell_t(\mathbf{w}_t))$
 - 9: **end for**
-

losses suffered by the algorithm and that of any fixed predictor, formally:

$$\text{Regret}_T(\mathbf{x}) := \sum_{t=1}^T f_t(\mathbf{x}_t) - \sum_{t=1}^T f_t(\mathbf{x}). \quad (1)$$

Yet, online learning sees the whole process as a single task and train only a single model, and doesn't consider the case where a new task is revealed in each iteration.

Neither paradigm alone is ideal for the continual lifelong learning scenario, thus, Finn et al. [FRKL19] proposed to combine them to construct the Online Meta-Learning framework which is presented in Algorithm 1. A meta-learner \mathbf{w}_t is maintained to represent the prior knowledge learned from past rounds. Upon seeing a new task \mathcal{T}_t , one is first given some training data \mathcal{D}_t^{tr} for adapting the meta-learner to the current task following some strategy $U(\cdot)$. Then the test data \mathcal{D}_t^{ts} will be revealed for evaluating the performance of the adapted learner $\hat{\mathbf{w}}_t$. The loss suffered at this round $\ell_t(\mathbf{w}_t)$ can then be fed into an online learning algorithm \mathcal{A} to update the meta-learner. We follow [FRKL19] and use $U(\mathbf{w}_t, \mathcal{D}_t^{tr}) = \mathbf{w}_t - \alpha \nabla \mathbb{E}_{\mathbf{x}, y \sim \mathcal{D}_t^{tr}} [\ell(\mathbf{w}_t, \mathbf{x}; y)]$ where α is the step-size.

As tasks can be very different, the original regret in Equation (1) of competing with a fixed learner across all tasks becomes less meaningful. Thus, Finn et al. [FRKL19] proposed a new regret as:

$$\text{Regret}'_T(\mathbf{w}) = \sum_{t=1}^T \ell(U(\mathbf{w}_t, \mathcal{D}_t^{tr}), \mathcal{D}_t^{ts}) - \sum_{t=1}^T \ell(U(\mathbf{w}, \mathcal{D}_t^{tr}), \mathcal{D}_t^{ts}),$$

which competes with any fixed *meta-learner*. Under this, they designed the Follow the Meta Leader algorithm enjoying a logarithmic regret when assuming strong-convexity on ℓ .

However, many problems of current interest have a non-convex nature. Thus, in Section 2, we generalize the online meta-learning framework to non-convex settings. We also explain why the regret of form (1) is infeasible in such case, and introduce an alternative performance measure. In Section 3 we exemplify the power of Algorithm 1 by using the AdaGrad-Norm [WWB19] algorithm as the online learning algorithm \mathcal{A} , and prove a logarithmic local regret robust to any hyperparameter initialization. Section 4 shows the empirical comparison of our algorithm with traditional methods.

Notation. We use bold letters to denote vectors, e.g., $\mathbf{u}, \mathbf{G} \in \mathbb{R}^d$. The i th coordinate of a vector \mathbf{u} is u_i . Unless explicitly noted, we study the Euclidean space \mathbb{R}^d with the inner product $\langle \cdot, \cdot \rangle$, and the Euclidean norm. We assume everywhere that our objective function f is bounded from below and denote the infimum by $f^* > -\infty$. The gradient of a function f at \mathbf{x} is denoted by $\nabla f(\mathbf{x})$. $\mathbb{E}[\mathbf{u}]$ means the expectation w.r.t. the underlying probability distribution of a random variable \mathbf{u} .

2 Generalizing Online Meta-Learning to Non-convex Regime

Finding the global minimum for a non-convex function in general is known to be NP-hard. Yet, if we could find an online learning algorithm with a $o(T)$ regret for some non-convex function classes, we can optimize any function f of that class efficiently: simply run the online learning algorithm but with the objective f as the loss ℓ_t at each round, and choose a random update as output. This gives us:

$$\mathbb{E}_i[f(\mathbf{w}_i)] - \min_{\mathbf{w} \in \mathcal{K}} f(\mathbf{w}) = \frac{1}{T} \sum_{t=1}^T f(\mathbf{w}_t) - \min_{\mathbf{w} \in \mathcal{K}} f(\mathbf{w}) = \frac{1}{T} \sum_{t=1}^T \ell_t(\mathbf{w}_t) - \min_{\mathbf{w} \in \mathcal{K}} \frac{1}{T} \sum_{t=1}^T \ell_t(\mathbf{w}) \in o(1).$$

This leads to a contradiction unless $P=NP$, and we are forced to find another performance measure for the non-convex case. One potential candidate is the local regret proposed by Hazan et al. [HSZ17]:

$$\mathcal{R}_m(T) \triangleq \sum_{t=1}^T \|\nabla F_{t,m}(\mathbf{w}_t)\|^2, \quad \text{where} \quad F_{t,m}(\mathbf{w}_t) \triangleq \frac{1}{m} \sum_{i=0}^{m-1} \ell_{t-i}(\mathbf{w}_t), \quad (2)$$

where $1 \leq m \leq T$, and $\ell_i(\cdot) = 0$ for $i \leq 0$.

The reason for using a sliding-window in F can be justified by Theorem 2.7 in [HSZ17].

3 Algorithm Design and Theoretical Guarantees

3.1 Stochasticity of Online Meta-learning Algorithms

In practice, \mathcal{D}_t^{ts} is typically just a random sample batch of the whole test-set, the losses and gradients obtained at each round are thus (unbiased) estimations of the true ones. This is the stochastic setting which we formalize by making following assumptions on obtaining the stochastic gradients.

Algorithm 2 AdaGrad-Norm

- 1: **Input:** Initialize $\mathbf{w}_1 \in \mathbb{R}^d$, $b_1 > 0$, $\eta > 0$.
 - 2: **for** $t = 1, \dots, T$ **do**
 - 3: Generate $\mathbf{G}_{t,m}(\mathbf{w}_t) = \frac{1}{m} \sum_{i=0}^{m-1} \mathbf{g}_{t-i}(\mathbf{w}_t, \xi_{t,t-i})$
 - 4: $b_{t+1}^2 \leftarrow b_t^2 + \|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2$
 - 5: $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \frac{\eta}{b_{t+1}} \mathbf{G}_{t,m}(\mathbf{w}_t)$
 - 6: **end for**
-

Assumption 1. We assume that at each round t , each call to any stochastic gradient oracle \mathbf{g}_i , $i \in \{t-m+1, \dots, t\}$, yields an i.i.d. random vector $\mathbf{g}_i(\mathbf{w}_t, \xi_{t,i})$ with the following properties:

- (a) *Unbiased:* $\mathbb{E}_{\xi_{t,i}}[\mathbf{g}_i(\mathbf{w}_t, \xi_{t,i}) | \xi_{1:t-1}] = \nabla \ell_i(\mathbf{w}_t)$;
- (b) *Bounded variance:* $\mathbb{E}_{\xi_{t,i}}[\|\mathbf{g}_i(\mathbf{w}_t, \xi_{t,i}) - \nabla \ell_i(\mathbf{w}_t)\|^2 | \xi_{1:t-1}] \leq \sigma^2$;
- (c) *Mutual independence:* for $i \neq j$,
 $\mathbb{E}_{\xi_{t,i}, \xi_{t,j}}[\langle \mathbf{g}_i(\mathbf{w}_t, \xi_{t,i}), \mathbf{g}_j(\mathbf{w}_t, \xi_{t,j}) \rangle | \xi_{1:t-1}] = \langle \mathbb{E}_{\xi_{t,i}}[\mathbf{g}_i(\mathbf{w}_t, \xi_{t,i}) | \xi_{1:t-1}], \mathbb{E}_{\xi_{t,j}}[\mathbf{g}_j(\mathbf{w}_t, \xi_{t,j}) | \xi_{1:t-1}] \rangle$.

where $\xi_{1:t-1} = \{\xi_{1,1}, \xi_{2,1}, \xi_{2,2}, \dots, \xi_{t-1,t-m}, \dots, \xi_{t-1,t-1}\}$, and $\mathbb{E}_{\xi_{t,i}}[\mathbf{u} | \xi_{1:t-1}]$ denotes the conditional expectation of \mathbf{u} with respect to $\xi_{1:t-1}$. Also note that $\mathbf{g}_i(\cdot) = 0$ for $i \leq 0$.

Hazan et al. proposed a time-smoothed online gradient descent algorithm [HSZ17] for such case. Yet, that algorithm's performance critically relies on the choice of the step-size η , and may even diverge if $\eta > \frac{2}{\beta}$ where β is the (often unknown) smoothness of the loss function. We thus propose to use the AdaGrad-Norm [WWB19] algorithm (Algorithm 2) as the online learning algorithm \mathcal{A} in Algorithm 1 instead. Here, $b_1 > 0$ is the initialization of the accumulated squared norms and prevents division by 0, while $\eta > 0$ is to ensure homogeneity and that the units match.

We present below an analysis of this algorithm assuming the loss function $\ell : \mathcal{K} \rightarrow \mathbb{R}$ satisfies:

Assumption 2. (C^2 -smoothness and boundedness)

- (a) *L-Lipschitzness:* $\forall \mathbf{u}, \mathbf{v} \in \mathcal{K}, \|\ell(\mathbf{u}) - \ell(\mathbf{v})\| \leq L$.
- (b) *β -smoothness:* ℓ is differentiable and $\forall \mathbf{u}, \mathbf{v} \in \mathcal{K}, \|\nabla \ell(\mathbf{u}) - \nabla \ell(\mathbf{v})\| \leq \beta \|\mathbf{u} - \mathbf{v}\|$.
 Note that this implies [Nes03, Lemma 1.2.3], $\forall \mathbf{u}, \mathbf{v} \in \mathcal{K}$

$$|\ell(\mathbf{v}) - \ell(\mathbf{u}) - \langle \nabla \ell(\mathbf{u}), \mathbf{v} - \mathbf{u} \rangle| \leq \frac{\beta}{2} \|\mathbf{v} - \mathbf{u}\|^2. \quad (3)$$
- (c) *H-Hessian-Lipschitzness:* $\forall \mathbf{u}, \mathbf{v} \in \mathcal{K}, \|\nabla^2 \ell(\mathbf{u}) - \nabla^2 \ell(\mathbf{v})\| \leq H \|\mathbf{u} - \mathbf{v}\|$.
- (d) *M-Boundedness:* $\forall \mathbf{u} \in \mathcal{K}, |\ell(\mathbf{u})| \leq M$

3.2 Theoretical Results

With these common assumptions, we are ready to show the convergence analysis of the algorithm. The following lemma converts Assumption 2 of ℓ to properties of ℓ_t (see proof in Appendices):

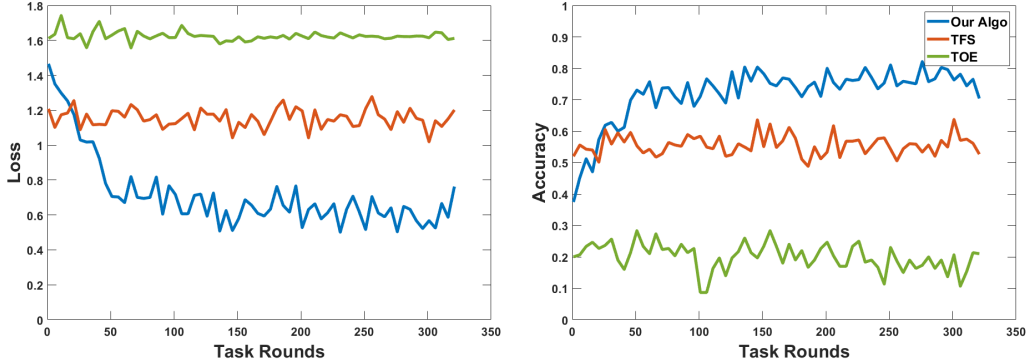


Figure 1: The comparison between our algorithm, TOE, and TFS on the Omniglot dataset.

Lemma 3. Assuming Assumption 2, ℓ_t is M -Bounded, $(1 + \alpha\beta)L$ -Lipschitz, and $(\alpha LH + (1 + \alpha\beta)^2\beta)$ -smooth.

For such loss functions, We have the following guarantee on our algorithm’s performance:

Theorem 4. Let ℓ_1, \dots, ℓ_T satisfy Assumptions 2. Then, feeding Algorithm 2 into Algorithm 1 with access to stochastic gradient oracles satisfying Assumptions 1 gives the following upper bound of $\mathcal{R}_m(T)$, with probability $1 - \delta$:

$$\mathcal{R}_m(T) \leq \frac{48C^2}{\delta^2} + \frac{8b_1C}{\delta} + \frac{8\sigma C\sqrt{T}}{\delta^{3/2}\sqrt{m}}.$$

where $C = \frac{4MT}{\eta m} + \left(\frac{\eta(\alpha LH + (1 + \alpha\beta)^2\beta) + 4\sigma/\sqrt{m}}{2} \right) \ln \left(1 + \frac{2(\sigma^2/m + (1 + \alpha\beta)^2 L^2)T}{b_1^2} \right)$.

By selecting $m \in \Theta(T)$, a logarithmic regret of the algorithm is guaranteed w.r.t. any $b_1, \eta > 0$.

4 Experiments

For evaluation, we applied the proposed algorithm to the few-shot image recognition task on the Omniglot [LST15] dataset. We employed the N -way K -shot protocol following [VBL⁺16]: at each round, pick N unseen characters irrespective of alphabets. Provide the meta-learner w_t with K different drawings of each of the N characters as the training set \mathcal{D}^{tr} , then evaluate the adapted model \hat{w}_t ’s ability on new unseen instances within the N classes (namely the test set \mathcal{D}^{ts}). We chose the 5-way 5-shot scheme, and used 15 examples per character for testing following [RL17]. We employed a CNN model from [VBL⁺16] whose detailed architecture can be found in Appendix A.3.

To study if our algorithm provides any empirical benefit over traditional methods, we compare it to two algorithms [FRKL19]: Train on Everything (TOE), and Train from Scratch (TFS). On each round t , both initialize a new model. The difference is that TOE trains over all available data, both training and testing, from all past tasks, plus \mathcal{D}_t^{tr} at current round, while TFS only uses \mathcal{D}_t^{tr} for training.

The result is shown in Figure 1 which suggests that our algorithm gradually accumulates prior knowledge, which enables fast learning of later tasks. TFS provides a good example of how CNN performs when the training data is scarce. On the contrary, TOE behaves nearly as random guessing. The inferiority of TOE to TFS is somehow surprising, as TOE has much more training data than TFS. The reason is that TOE regards all training data as coming from a single distribution, and tries to learn a model that works for all tasks. Thus, when tasks are substantially different from each other, TOE might even incur negative transfer and fail to solve any single task as has been observed in [PBS16]. Meanwhile, by using training data of the current task only, TFS avoids negative transfer, but also rules out learning of any connection between tasks. Our algorithm, in contrast, is designed to discover common structures across tasks, and use these information to guide fast adaptation to new tasks.

5 Conclusion

The continual lifelong learning problem is common in real-life, where an agent needs to accumulate knowledge from every task it encounters, and utilizes that knowledge for fast learning of new tasks. To solve this problem, we can combine the meta-learning and the online-learning paradigms to form the online meta-learning framework. In this work, we generalized this framework to the non-convex setting, and introduced the local regret to replace the original regret definition. We applied it to the stochastic setting, and showed its superiority both in theory and practice.

References

- [AAA⁺16] Dario Amodei, Sundaram Ananthanarayanan, Rishita Anubhai, Jingliang Bai, Eric Battenberg, Carl Case, Jared Casper, Bryan Catanzaro, Qiang Cheng, Guoliang Chen, et al. Deep speech 2: End-to-end speech recognition in english and mandarin. In *Proceedings of International Conference on Machine Learning*, pages 173–182, 2016.
- [CBL06] Nicolo Cesa-Bianchi and Gabor Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.
- [FRKL19] Chelsea Finn, Aravind Rajeswaran, Sham Kakade, and Sergey Levine. Online meta-learning. In *Proceedings of International Conference on Machine Learning*, pages 1920–1930, 2019.
- [HSZ17] Elad E Hazan, Karan Singh, and Cyril Zhang. Efficient regret minimization in non-convex games. In *Proceedings of International Conference on Machine Learning*, pages 2278–2288, 2017.
- [IS15] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of International Conference on Machine Learning*, pages 448–456, 2015.
- [KB14] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2014.
- [LBH15] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436, 2015.
- [LO19] Xiaoyu Li and Francesco Orabona. On the convergence of stochastic gradient descent with adaptive stepsizes. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 983–992, 2019.
- [LST15] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- [LUTG17] Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. Building machines that learn and think like people. *Behavioral and Brain Sciences*, 40, 2017.
- [Nes03] Y. Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer, 2003.
- [NM92] Devang K Naik and RJ Mammone. Meta-neural networks that learn by learning. In *Proceedings of International Joint Conference on Neural Networks*, volume 1, pages 437–442. IEEE, 1992.
- [PBS16] Emilio Parisotto, Jimmy Lei Ba, and Ruslan Salakhutdinov. Actor-mimic: Deep multitask and transfer reinforcement learning. *International Conference on Learning Representations*, 2016.
- [PGC⁺17] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- [RDGF16] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 779–788, 2016.
- [RL17] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. *International Conference on Learning Representations*, 2017.
- [SBB⁺16] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-learning with memory-augmented neural networks. In *Proceedings of International Conference on Machine Learning*, pages 1842–1850, 2016.

- [SHM⁺16] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484, 2016.
- [TP12] Sebastian Thrun and Lorien Pratt. *Learning to learn*. Springer Science & Business Media, 2012.
- [VBL⁺16] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *Proceedings of Advances in Neural Information Processing Systems*, pages 3630–3638, 2016.
- [WWB19] Rachel Ward, Xiaoxia Wu, and Leon Bottou. Adagrad stepsizes: sharp convergence over nonconvex landscapes. In *Proceedings of International Conference on Machine Learning*, pages 6677–6686, 2019.

A Appendix

A.1 Proof of Lemma 3

Proof. We first write out the complete formula of ℓ_t :

$$\begin{aligned}\ell_t(\mathbf{w}) &= \ell(\hat{\mathbf{w}}, \mathcal{D}_t^{ts}) \\ &= \mathbb{E}_{\mathbf{x}^{ts}, y^{ts} \sim \mathcal{D}_t^{ts}} [\ell(\mathbf{U}(\mathbf{w}, \mathcal{D}_t^{tr}), \mathbf{x}^{ts}; y^{ts})] \\ &= \mathbb{E}_{\mathbf{x}^{ts}, y^{ts} \sim \mathcal{D}_t^{ts}} [\ell(\mathbf{w} - \alpha \nabla \mathbb{E}_{\mathbf{x}^{tr}, y^{tr} \sim \mathcal{D}_t^{tr}} [\ell(\mathbf{w}, \mathbf{x}^{tr}; y^{tr})], \mathbf{x}^{ts}; y^{ts})] \\ &\triangleq f_t(\mathbf{w} - \alpha \nabla \hat{f}_t(\mathbf{w})) .\end{aligned}$$

The M -Boundedness is straight-forward.

To show the Lipschitzness, we derive $\nabla \ell_t$:

$$\nabla \ell_t(\mathbf{w}) = (\mathbf{I} - \alpha \nabla^2 \hat{f}_t(\mathbf{w})) \nabla f_t(\mathbf{w} - \alpha \nabla \hat{f}_t(\mathbf{w})) .$$

Note that f_t and \hat{f}_t both share the properties of ℓ , thus, from Assumption 2(a,b), we have:

$$\|\nabla \ell_t(\mathbf{w})\| \leq (1 + \alpha\beta) \|\nabla f_t(\mathbf{w} - \alpha \nabla \hat{f}_t(\mathbf{w}))\| \leq (1 + \alpha\beta)L .$$

Next, denoting $\mathbf{U}(\mathbf{w}, \mathcal{D}_t^{tr})$ as $\mathbf{U}_t(\mathbf{w})$, we have $\forall \mathbf{u}, \mathbf{v} \in \mathcal{K}$:

$$\begin{aligned}& \|\nabla \ell_t(\mathbf{u}) - \nabla \ell_t(\mathbf{v})\| \\ &= \|\nabla \mathbf{U}_t(\mathbf{u}) \nabla f_t(\mathbf{U}_t(\mathbf{u})) - \nabla \mathbf{U}_t(\mathbf{v}) \nabla f_t(\mathbf{U}_t(\mathbf{v}))\| \\ &= \|\nabla \mathbf{U}_t(\mathbf{u}) \nabla f_t(\mathbf{U}_t(\mathbf{u})) - \nabla \mathbf{U}_t(\mathbf{v}) \nabla f_t(\mathbf{U}_t(\mathbf{u})) + \nabla \mathbf{U}_t(\mathbf{v}) \nabla f_t(\mathbf{U}_t(\mathbf{u})) - \nabla \mathbf{U}_t(\mathbf{v}) \nabla f_t(\mathbf{U}_t(\mathbf{v}))\| \\ &\leq \|(\nabla \mathbf{U}_t(\mathbf{u}) - \nabla \mathbf{U}_t(\mathbf{v})) \nabla f_t(\mathbf{U}_t(\mathbf{u}))\| + \|\nabla \mathbf{U}_t(\mathbf{v}) (\nabla f_t(\mathbf{U}_t(\mathbf{u})) - \nabla f_t(\mathbf{U}_t(\mathbf{v})))\| \\ &= \alpha \|(\nabla^2 \hat{f}_t(\mathbf{u}) - \nabla^2 \hat{f}_t(\mathbf{v})) \nabla f_t(\mathbf{U}_t(\mathbf{u}))\| + \|(\mathbf{I} - \alpha \nabla^2 \hat{f}_t(\mathbf{v})) (\nabla f_t(\mathbf{U}_t(\mathbf{u})) - \nabla f_t(\mathbf{U}_t(\mathbf{v})))\| \\ &\leq \alpha LH \|\mathbf{u} - \mathbf{v}\| + (1 + \alpha\beta) \|\nabla f_t(\mathbf{U}_t(\mathbf{u})) - \nabla f_t(\mathbf{U}_t(\mathbf{v}))\| \\ &\leq \alpha LH \|\mathbf{u} - \mathbf{v}\| + (1 + \alpha\beta)\beta \|\mathbf{U}_t(\mathbf{u}) - \mathbf{U}_t(\mathbf{v})\| \\ &\leq \alpha LH \|\mathbf{u} - \mathbf{v}\| + (1 + \alpha\beta)^2 \beta \|\mathbf{u} - \mathbf{v}\| .\end{aligned}$$

where the first inequality uses the triangle inequality of a norm; the second inequality uses the smoothness and hessian-Lipschitzness assumptions; the third inequality uses the smoothness assumption.

We are left to prove the last inequality:

$$\begin{aligned}\|\mathbf{U}_t(\mathbf{u}) - \mathbf{U}_t(\mathbf{v})\| &= \|\mathbf{u} - \alpha \nabla \hat{f}_t(\mathbf{u}) - \mathbf{v} + \alpha \nabla \hat{f}_t(\mathbf{v})\| \\ &= \|\mathbf{u} - \mathbf{v} - \alpha (\nabla \hat{f}_t(\mathbf{u}) - \nabla \hat{f}_t(\mathbf{v}))\| \\ &\leq \|\mathbf{u} - \mathbf{v}\| + \alpha \|\nabla \hat{f}_t(\mathbf{u}) - \nabla \hat{f}_t(\mathbf{v})\| \\ &\leq (1 + \alpha\beta) \|\mathbf{u} - \mathbf{v}\| .\end{aligned}$$

where the the first inequality uses the triangle inequality of a norm, and the second inequality uses the smoothness assumption. \square

A.2 Proof of Theorem 4

For simplicity, we denote \mathbb{E}_t as condition on $\xi_{1:t-1}$ and take expectation w.r.t. $\xi_{t,t-m+1}, \dots, \xi_{t,t}$, where $\xi_{1:t-1} = \{\xi_{1,1}, \xi_{2,1}, \xi_{2,2}, \dots, \xi_{t-1,t-m}, \dots, \xi_{t-1,t-1}\}$.

For the proof, we need the following technical lemmas.

Lemma 5. Recall $\mathbf{G}_{t,m}(\mathbf{w}_t) = \frac{1}{m} \sum_{i=0}^{m-1} \mathbf{g}_{t-i}(\mathbf{w}_t, \xi_{t,t-i})$, and $F_{t,m}(\mathbf{w}_t) = \frac{1}{m} \sum_{i=0}^{m-1} \ell_{t-i}(\mathbf{w}_t)$, then Assumption 1 gives us:

$$(a) \text{ Unbiased: } \mathbb{E}_t [\mathbf{G}_{t,m}(\mathbf{w}_t)] = \nabla F_{t,m}(\mathbf{w}_t)$$

$$(b) \text{ Bounded variance: } \mathbb{E}_t \left[\|\mathbf{G}_{t,m}(\mathbf{x}_t) - \nabla F_{t,m}(\mathbf{x}_t)\|^2 \right] \leq \frac{\sigma^2}{m}$$

Proof. In Assumption 1(a) we assume $\mathbb{E}_{\xi_{t,i}} [\mathbf{g}_i(\mathbf{w}_t, \xi_{t,i}) | \xi_{1:t-1}] = \nabla \ell_i(\mathbf{w}_t)$ for $i \in \{t-m+1, \dots, t\}$, the linearity of expectation immediately gives us $\mathbb{E}_t [\mathbf{G}_{t,m}(\mathbf{w}_t)] = \nabla F_{t,m}(\mathbf{w}_t)$.

To see the second part, we only need to expand $\mathbb{E}_t [\|\mathbf{G}_{t,m}(\mathbf{x}_t) - \nabla F_{t,m}(\mathbf{x}_t)\|^2]$ as:

$$\begin{aligned} & \frac{1}{m^2} \mathbb{E}_t \left[\left\| \sum_{i=0}^{m-1} \mathbf{g}_{t-i}(\mathbf{w}_t, \xi_{t,t-i}) - \nabla \ell_{t-i}(\mathbf{w}_t) \right\|^2 \right] \\ &= \frac{1}{m^2} \sum_{i=0}^{m-1} \sum_{j=0}^{m-1} \mathbb{E}_t [\langle \mathbf{g}_{t-i}(\mathbf{w}_t, \xi_{t,t-i}) - \nabla \ell_{t-i}(\mathbf{w}_t), \mathbf{g}_{t-j}(\mathbf{w}_t, \xi_{t,t-j}) - \nabla \ell_{t-j}(\mathbf{w}_t) \rangle] \\ &= \frac{1}{m^2} \sum_{i=0}^{m-1} \mathbb{E}_t [\|\mathbf{g}_{t-i}(\mathbf{w}_t, \xi_{t,t-i}) - \nabla \ell_{t-i}(\mathbf{w}_t)\|^2] \\ & \quad + \frac{1}{m^2} \sum_{i=0}^{m-1} \sum_{j \neq i}^{m-1} \mathbb{E}_t [\langle \mathbf{g}_{t-i}(\mathbf{w}_t, \xi_{t,t-i}) - \nabla \ell_{t-i}(\mathbf{w}_t), \mathbf{g}_{t-j}(\mathbf{w}_t, \xi_{t,t-j}) - \nabla \ell_{t-j}(\mathbf{w}_t) \rangle] . \end{aligned}$$

Each item of the first part in the last equation can be bounded by σ^2 according to Assumption 1(b), which leads to a $\frac{\sigma^2}{m}$ overall upper-bound.

For the second part, we need to use the Mutual Independence assumption (namely Assumption 1(c)):

$$\begin{aligned} & \mathbb{E}_t [\langle \mathbf{g}_{t-i}(\mathbf{w}_t, \xi_{t,t-i}) - \nabla \ell_{t-i}(\mathbf{w}_t), \mathbf{g}_{t-j}(\mathbf{w}_t, \xi_{t,t-j}) - \nabla \ell_{t-j}(\mathbf{w}_t) \rangle] \\ &= \langle \mathbb{E}_t [\mathbf{g}_{t-i}(\mathbf{w}_t, \xi_{t,t-i}) - \nabla \ell_{t-i}(\mathbf{w}_t)], \mathbb{E}_t [\mathbf{g}_{t-j}(\mathbf{w}_t, \xi_{t,t-j}) - \nabla \ell_{t-j}(\mathbf{w}_t)] \rangle . \end{aligned}$$

Use Assumption 1(a) again we see that the above equation equals 0.

This proves part (b) of this lemma. \square

Lemma 6. *Given Assumption 2(d), we have: $\sum_{t=1}^T \mathbb{E}[F_{t,m}(\mathbf{w}_t) - F_{t,m}(\mathbf{w}_{t+1})] \leq \frac{4MT}{m}$.*

Proof.

$$\begin{aligned} & \sum_{t=1}^T \mathbb{E}[F_{t,m}(\mathbf{w}_t) - F_{t,m}(\mathbf{w}_{t+1})] \\ &= \sum_{t=2}^T \mathbb{E}[F_{t,m}(\mathbf{w}_t) - F_{t-1,m}(\mathbf{w}_t)] + F_{1,m}(\mathbf{w}_1) - \mathbb{E}[F_{T,m}(\mathbf{w}_{T+1})] \\ &= \sum_{t=2}^T \frac{1}{m} \sum_{i=0}^{m-1} \mathbb{E}[\ell_{t-i}(\mathbf{w}_t) - \ell_{t-1-i}(\mathbf{w}_t)] + \ell_1(\mathbf{w}_1) - \frac{1}{m} \sum_{i=0}^{m-1} \mathbb{E}[\ell_{T-i}(\mathbf{w}_{T+1})] \\ &= \sum_{t=2}^T \frac{1}{m} \mathbb{E}[\ell_t(\mathbf{w}_t) - \ell_{t-m}(\mathbf{w}_t)] + \ell_1(\mathbf{w}_1) - \frac{1}{m} \sum_{i=0}^{m-1} \mathbb{E}[\ell_{T-i}(\mathbf{w}_{T+1})] \\ &\leq \frac{2MT}{m} + M + M \leq \frac{4MT}{m} , \end{aligned}$$

where we use the definition that $\ell_i(\cdot) = 0$ for $i \leq 0$, and $1 \leq m \leq T$. \square

Lemma 7 ([LO19], Lemma 9). *Let $h : [0, +\infty) \rightarrow [0, +\infty)$ be a nonincreasing function, and $a_i \geq 0$ for $i = 0, \dots, T$. Then*

$$\sum_{t=1}^T a_t h \left(a_0 + \sum_{i=1}^t a_i \right) \leq \int_{a_0}^{\sum_{t=0}^T a_t} h(x) dx .$$

Proof. Denote $s_t = \sum_{i=0}^t a_i$.

$$a_t h(s_t) = \int_{s_{t-1}}^{s_t} h(s_t) dx \leq \int_{s_{t-1}}^{s_t} h(x) dx .$$

Summing over $t = 1, \dots, T$, we have the stated bound. \square

Proof of Theorem 4. The proof presented below closely follows that of Theorem 2.1 in [WWB19].

For simplicity, we denote each ℓ_t L' -Lipschitz, β' -smooth, and M -bounded.

Given that $F_{t,m}$ is the average of m β' -smooth functions, it can be easily shown that $F_{t,m}$ is also β' -smooth. Thus, from the property of a smooth function (Equation (3)), we have:

$$\begin{aligned} \frac{F_{t,m}(\mathbf{w}_{t+1}) - F_{t,m}(\mathbf{w}_t)}{\eta} &\leq \frac{1}{\eta} \left[\langle \nabla F_{t,m}(\mathbf{w}_t), \mathbf{w}_{t+1} - \mathbf{w}_t \rangle + \frac{\beta'}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 \right] \\ &= -\langle \nabla F_{t,m}(\mathbf{w}_t), \frac{\mathbf{G}_{t,m}(\mathbf{w}_t)}{b_{t+1}} \rangle + \frac{\eta\beta'}{2b_{t+1}^2} \|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2. \end{aligned}$$

Take expectation w.r.t. $\xi_{t,t-m+1}, \dots, \xi_{t,t}$ conditioned on $\xi_{1:t-1}$ (namely $\mathbb{E}_t[\cdot]$) on both sides:

$$\begin{aligned} &\frac{\mathbb{E}_t[F_{t,m}(\mathbf{w}_{t+1}) - F_{t,m}(\mathbf{w}_t)]}{\eta} = \frac{\mathbb{E}_t[F_{t,m}(\mathbf{w}_{t+1})] - F_{t,m}(\mathbf{w}_t)}{\eta} \\ &\leq \mathbb{E}_t \left[-\langle \nabla F_{t,m}(\mathbf{w}_t), \frac{\mathbf{G}_{t,m}(\mathbf{w}_t)}{b_{t+1}} \rangle \right] + \frac{\eta\beta'}{2} \mathbb{E}_t \left[\frac{1}{b_{t+1}^2} \|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2 \right] \\ &= \mathbb{E}_t \left[\frac{\langle \nabla F_{t,m}(\mathbf{w}_t), \mathbf{G}_{t,m}(\mathbf{w}_t) \rangle}{\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} - \frac{\langle \nabla F_{t,m}(\mathbf{w}_t), \mathbf{G}_{t,m}(\mathbf{w}_t) \rangle}{b_{t+1}} \right] \quad (1) \\ &\quad - \frac{\|\nabla F_{t,m}(\mathbf{w}_t)\|^2}{\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} + \frac{\eta\beta'}{2} \mathbb{E}_t \left[\frac{1}{b_{t+1}^2} \|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2 \right]. \quad (2) \end{aligned}$$

where the last equality is based on the unbiasedness assumption of $\mathbf{G}_{t,m}(\mathbf{w}_t)$ in Lemma 5 (a).

Now,

$$\begin{aligned} &\left| \frac{1}{\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} - \frac{1}{b_{t+1}} \right| \\ &= \left| \frac{b_{t+1} - \sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}}{b_{t+1}\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} \right| \\ &= \left| \frac{b_{t+1}^2 - b_t^2 - \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 - \sigma^2/m}{b_{t+1}\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m} (b_{t+1} + \sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m})} \right| \\ &= \left| \frac{\|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2 - \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 - \sigma^2/m}{b_{t+1}\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m} (b_{t+1} + \sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m})} \right| \\ &= \frac{|\|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2 - \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 - \sigma^2/m|}{b_{t+1}\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m} (b_{t+1} + \sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m})} \\ &\leq \frac{|\|\mathbf{G}_{t,m}(\mathbf{w}_t)\| - \|\nabla F_{t,m}(\mathbf{w}_t)\|| (\|\mathbf{G}_{t,m}(\mathbf{w}_t)\| + \|\nabla F_{t,m}(\mathbf{w}_t)\|) + \sigma^2/m}{b_{t+1}\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m} (b_{t+1} + \sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m})} \\ &\leq \frac{|\|\mathbf{G}_{t,m}(\mathbf{w}_t)\| - \|\nabla F_{t,m}(\mathbf{w}_t)\||}{b_{t+1}\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} + \frac{\sigma/\sqrt{m}}{b_{t+1}\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}}, \end{aligned}$$

where the last inequality uses the fact that $\|\mathbf{G}_{t,m}(\mathbf{w}_t)\| \leq b_{t+1}$.

As $|\cdot|$ is a convex function, using Jensen's inequality, we can bound Equation (1) as:

$$\begin{aligned}
& \mathbb{E}_t \left[\frac{\langle \nabla F_{t,m}(\mathbf{w}_t), \mathbf{G}_{t,m}(\mathbf{w}_t) \rangle}{\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} - \frac{\langle \nabla F_{t,m}(\mathbf{w}_t), \mathbf{G}_{t,m}(\mathbf{w}_t) \rangle}{b_{t+1}} \right] \\
& \leq \left| \mathbb{E}_t \left[\left(\frac{1}{\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} - \frac{1}{b_{t+1}} \right) \langle \nabla F_{t,m}(\mathbf{w}_t), \mathbf{G}_{t,m}(\mathbf{w}_t) \rangle \right] \right| \\
& \leq \mathbb{E}_t \left[\left| \left(\frac{1}{\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} - \frac{1}{b_{t+1}} \right) \langle \nabla F_{t,m}(\mathbf{w}_t), \mathbf{G}_{t,m}(\mathbf{w}_t) \rangle \right| \right] \\
& = \mathbb{E}_t \left[\left| \left(\frac{1}{\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} - \frac{1}{b_{t+1}} \right) |\langle \nabla F_{t,m}(\mathbf{w}_t), \mathbf{G}_{t,m}(\mathbf{w}_t) \rangle| \right| \right] \\
& \leq \frac{\mathbb{E}_t [\|\mathbf{G}_{t,m}(\mathbf{w}_t)\| - \|\nabla F_{t,m}(\mathbf{w}_t)\|] \|\mathbf{G}_{t,m}(\mathbf{w}_t)\| \|\nabla F_{t,m}(\mathbf{w}_t)\|}{b_{t+1} \sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} \tag{3}
\end{aligned}$$

$$+ \frac{\mathbb{E}_t [\|\mathbf{G}_{t,m}(\mathbf{w}_t)\| \|\nabla F_{t,m}(\mathbf{w}_t)\| \sigma / \sqrt{m}]}{b_{t+1} \sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}}, \tag{4}$$

where in the last inequality we used the Cauchy-Schwarz inequality.

Next, applying the inequality $ab \leq \frac{\lambda}{2}a^2 + \frac{1}{2\lambda}b^2$ with $\lambda = \frac{2\sigma^2/m}{\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}}$, $a = \frac{\|\mathbf{G}_{t,m}(\mathbf{w}_t)\|}{b_{t+1}}$, and $b = \frac{\|\mathbf{G}_{t,m}(\mathbf{w}_t)\| - \|\nabla F_{t,m}(\mathbf{w}_t)\| \|\nabla F_{t,m}(\mathbf{w}_t)\|}{\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}}$, Equation (3) becomes:

$$\begin{aligned}
& \frac{\mathbb{E}_t [\|\mathbf{G}_{t,m}(\mathbf{w}_t)\| - \|\nabla F_{t,m}(\mathbf{w}_t)\|] \|\mathbf{G}_{t,m}(\mathbf{w}_t)\| \|\nabla F_{t,m}(\mathbf{w}_t)\|}{b_{t+1} \sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} \\
& \leq \frac{\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}}{4\sigma^2/m} \frac{\|\nabla F_{t,m}(\mathbf{w}_t)\|^2 \mathbb{E}_t [\|\mathbf{G}_{t,m}(\mathbf{w}_t)\| - \|\nabla F_{t,m}(\mathbf{w}_t)\|]^2}{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m} \\
& \quad + \frac{\sigma^2/m}{\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} \mathbb{E}_t \left[\frac{\|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2}{b_{t+1}^2} \right] \\
& \leq \frac{\|\nabla F_{t,m}(\mathbf{w}_t)\|^2}{4\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} + \frac{\sigma}{\sqrt{m}} \mathbb{E}_t \left[\frac{\|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2}{b_{t+1}^2} \right],
\end{aligned}$$

where the last inequality is because $\|\mathbf{u}\| - \|\mathbf{v}\| \leq \|\mathbf{u} - \mathbf{v}\|$ holds for any $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$.

Similarly, applying the inequality $ab \leq \frac{\lambda}{2}a^2 + \frac{1}{2\lambda}b^2$ with $\lambda = \frac{2}{\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}}$, $a = \frac{\|\mathbf{G}_{t,m}(\mathbf{w}_t)\| \sigma / \sqrt{m}}{b_{t+1}}$, and $b = \frac{\|\nabla F_{t,m}(\mathbf{w}_t)\|}{\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}}$, Equation (4) becomes:

$$\begin{aligned}
& \frac{\mathbb{E}_t [\|\mathbf{G}_{t,m}(\mathbf{w}_t)\| \|\nabla F_{t,m}(\mathbf{w}_t)\| \sigma / \sqrt{m}]}{b_{t+1} \sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} \\
& \leq \frac{\|\nabla F_{t,m}(\mathbf{w}_t)\|^2}{4\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} + \frac{\sigma}{\sqrt{m}} \mathbb{E}_t \left[\frac{\|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2}{b_{t+1}^2} \right].
\end{aligned}$$

Putting above two inequalities back, we have:

$$\begin{aligned}
& \mathbb{E}_t \left[\frac{\langle \nabla F_{t,m}(\mathbf{w}_t), \mathbf{G}_{t,m}(\mathbf{w}_t) \rangle}{\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} - \frac{\langle \nabla F_{t,m}(\mathbf{w}_t), \mathbf{G}_{t,m}(\mathbf{w}_t) \rangle}{b_{t+1}} \right] \\
& \leq \frac{\|\nabla F_{t,m}(\mathbf{w}_t)\|^2}{2\sqrt{b_t^2 + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m}} + \frac{2\sigma}{\sqrt{m}} \mathbb{E}_t \left[\frac{\|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2}{b_{t+1}^2} \right].
\end{aligned}$$

And put this back into Equation (1):

$$\begin{aligned} \frac{\mathbb{E}_t[F_{t,m}(\mathbf{w}_{t+1})] - F_{t,m}(\mathbf{w}_t)}{\eta} &\leq -\frac{\|\nabla F_{t,m}(\mathbf{w}_t)\|^2}{2\sqrt{b_t^2} + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m} \\ &\quad + \left(\frac{\eta\beta'}{2} + \frac{2\sigma}{\sqrt{m}}\right) \mathbb{E}_t \left[\frac{1}{b_{t+1}^2} \|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2 \right]. \end{aligned}$$

which can be rearranged as:

$$\begin{aligned} \frac{\|\nabla F_{t,m}(\mathbf{w}_t)\|^2}{2\sqrt{b_t^2} + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m} &\leq \frac{F_{t,m}(\mathbf{w}_t) - \mathbb{E}_t[F_{t,m}(\mathbf{w}_{t+1})]}{\eta} \\ &\quad + \left(\frac{\eta\beta'}{2} + \frac{2\sigma}{\sqrt{m}}\right) \mathbb{E}_t \left[\frac{1}{b_{t+1}^2} \|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2 \right]. \end{aligned}$$

Using the law of total expectation, we can take expectation with respect to $\xi_{1:t-1}$ on both sides:

$$\begin{aligned} \mathbb{E} \left[\frac{\|\nabla F_{t,m}(\mathbf{w}_t)\|^2}{2\sqrt{b_t^2} + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m} \right] &\leq \frac{\mathbb{E}[F_{t,m}(\mathbf{w}_t)] - \mathbb{E}[F_{t,m}(\mathbf{w}_{t+1})]}{\eta} \\ &\quad + \left(\frac{\eta\beta'}{2} + \frac{2\sigma}{\sqrt{m}}\right) \mathbb{E} \left[\frac{1}{b_{t+1}^2} \|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2 \right]. \end{aligned}$$

Now summing both sides from $t = 1$ to T :

$$\sum_{t=1}^T \mathbb{E} \left[\frac{\|\nabla F_{t,m}(\mathbf{w}_t)\|^2}{2\sqrt{b_t^2} + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m} \right] \leq \frac{\sum_{t=1}^T [\mathbb{E}[F_{t,m}(\mathbf{w}_t)] - \mathbb{E}[F_{t,m}(\mathbf{w}_{t+1})]]}{\eta} \quad (5)$$

$$+ \left(\frac{\eta\beta' + 4\sigma/\sqrt{m}}{2}\right) \mathbb{E} \sum_{t=1}^T \frac{\|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2}{b_{t+1}^2}. \quad (6)$$

As $b_{t+1}^2 = b_1^2 + \sum_{i=1}^t \|\mathbf{G}_{i,m}(\mathbf{w}_i)\|^2$, letting $h(x)$ be $1/x$ in Lemma 7 gives us:

$$\sum_{t=1}^T \frac{\|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2}{b_{t+1}^2} \leq \int_{b_1^2}^{b_{T+1}^2} \frac{1}{x} dx = \ln \left(\frac{b_{T+1}^2}{b_1^2} \right) = \ln \left(1 + \frac{\sum_{t=1}^T \|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2}{b_1^2} \right).$$

As $\ln(x)$ is a concave function in $(0, +\infty)$, using Jensen's inequality gives us:

$$\mathbb{E} \ln \left(1 + \frac{\sum_{t=1}^T \|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2}{b_1^2} \right) \leq \ln \left(1 + \frac{\sum_{t=1}^T \mathbb{E} [\|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2]}{b_1^2} \right)$$

Since each ℓ_t is L' -Lipschitz, so is $F_{t,m}(\cdot)$, thus, using Cauchy-Schwartz inequality:

$$\begin{aligned} \mathbb{E} [\|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2] &= \mathbb{E} [\|\mathbf{G}_{t,m}(\mathbf{w}_t) - \nabla F_{t,m}(\mathbf{w}_t) + \nabla F_{t,m}(\mathbf{w}_t)\|^2] \\ &\leq 2\mathbb{E} [\|\mathbf{G}_{t,m}(\mathbf{w}_t) - \nabla F_{t,m}(\mathbf{w}_t)\|^2] + 2\mathbb{E} [\|\nabla F_{t,m}(\mathbf{w}_t)\|^2] \\ &\leq 2(\sigma^2/m + L'^2). \end{aligned}$$

Put the above inequality back into Equation (6) and Lemma 6 back into Equation (5), we have:

$$\sum_{t=1}^T \mathbb{E} \left[\frac{\|\nabla F_{t,m}(\mathbf{w}_t)\|^2}{2\sqrt{b_t^2} + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m} \right] \quad (7)$$

$$\leq \frac{4MT}{\eta m} + \left(\frac{\eta\beta' + 4\sigma/\sqrt{m}}{2}\right) \ln \left(1 + \frac{2(\sigma^2/m + L'^2)T}{b_1^2} \right). \quad (8)$$

Using Markov's inequality, we have that with probability $1 - \delta_1$:

$$\sum_{t=1}^T \|\nabla F_{t,m}(\mathbf{w}_t) - \mathbf{G}_{t,m}(\mathbf{w}_t)\|^2 \leq \frac{T\sigma^2}{m\delta_1}.$$

Denote $Z \triangleq \sum_{t=1}^T \|\nabla F_{t,m}(\mathbf{w}_t)\|^2$, then with probability $1 - \delta_1$:

$$\begin{aligned}
& b_T^2 + \|\nabla F_{T,m}(\mathbf{w}_T)\|^2 + \sigma^2/m \\
&= b_1^2 + \sum_{t=1}^{T-1} \|\mathbf{G}_{t,m}(\mathbf{w}_t)\|^2 + \|\nabla F_{T,m}(\mathbf{w}_T)\|^2 + \sigma^2/m \\
&\leq b_1^2 + 2 \sum_{t=1}^T \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + 2 \sum_{t=1}^{T-1} \|\nabla F_{t,m}(\mathbf{w}_t) - \mathbf{G}_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m \\
&\leq b_1^2 + 2Z + 2T \frac{\sigma^2}{m\delta_1}
\end{aligned}$$

This means, with probability $1 - \delta_1$, we have:

$$\begin{aligned}
\sum_{t=1}^T \frac{\|\nabla F_{t,m}(\mathbf{w}_t)\|^2}{2\sqrt{b_t^2} + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m} &\geq \frac{\sum_{t=1}^T \|\nabla F_{t,m}(\mathbf{w}_t)\|^2}{2\sqrt{b_T^2} + \|\nabla F_{T,m}(\mathbf{w}_T)\|^2 + Z + \sigma^2/m} \\
&\geq \frac{\sum_{t=1}^T \|\nabla F_{t,m}(\mathbf{w}_t)\|^2}{2\sqrt{b_1^2} + 3Z + 2T \frac{\sigma^2}{m\delta_1}}.
\end{aligned}$$

Denote Equation (8) as C , and use Markov's inequality again we have, with probability $1 - \delta_2$:

$$\sum_{t=1}^T \frac{\|\nabla F_{t,m}(\mathbf{w}_t)\|^2}{2\sqrt{b_t^2} + \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 + \sigma^2/m} \leq \frac{C}{\delta_2}.$$

Therefore, with probability $1 - \delta_1 - \delta_2$ (as $(1 - \delta_1)(1 - \delta_2) \geq 1 - \delta_1 - \delta_2$):

$$\frac{Z}{2\sqrt{b_1^2} + 3Z + 2T \frac{\sigma^2}{m\delta_1}} \leq \frac{C}{\delta_2}.$$

This is equivalent to:

$$Z^2 - \frac{12C^2}{\delta_2^2} Z - \frac{4C^2}{\delta_2^2} \left(b_1^2 + \frac{2\sigma^2 T}{m\delta_1} \right) \leq 0.$$

Solving the above quadratic inequality of Z we have:

$$\begin{aligned}
Z &\leq \frac{6C^2}{\delta_2^2} + \sqrt{\frac{36C^4}{\delta_2^4} + \frac{4C^2}{\delta_2^2} \left(b_1^2 + \frac{2\sigma^2 T}{m\delta_1} \right)} \\
&\leq \frac{12C^2}{\delta_2^2} + \frac{2C}{\delta_2} \left(b_1 + \frac{\sigma\sqrt{2T}}{\sqrt{m\delta_1}} \right).
\end{aligned}$$

Let $\delta_1 = \delta_2 = \frac{\delta}{2}$, we have that with probability $1 - \delta$:

$$\sum_{t=1}^T \|\nabla F_{t,m}(\mathbf{w}_t)\|^2 \leq \frac{48C^2}{\delta^2} + \frac{4C}{\delta} \left(b_1 + \frac{2\sigma\sqrt{T}}{\sqrt{m\delta}} \right).$$

□

A.3 Experiment details

The Omniglot dataset consists of 20 instances of 1623 characters from 50 different alphabets, where each instance was drawn by a different person. Following [SBB⁺16], we augmented the data set with rotations by multiples of 90 degrees.

The CNN model we employed comes from [VBL⁺16]. It contains 4 modules, each of which is a 3x3 convolution with 64 filters followed by batch normalization [IS15], a ReLu non-linearity and 2x2 max-pooling. Images are downsampled to 28x28 so that the resulting feature map of the last hidden layer is 1x1x64. The last layer is fed into a fully connected layer and the loss we used is the Cross-Entropy loss.

The experiments are done in PyTorch [PGC⁺17], and parameters are by default if no specification is provided. For the parameter α in the local adapter strategy $U(\cdot)$ in Algorithm 1, we set it to be 0.1 everywhere, and the gradient descent step is performed only once for each task. For the AdaGrad-Norm algorithm (Algorithm 2) we used, we set $b_1 = \eta = 1$ as suggested in the original paper [WWB19]. The TFS and TOE used Adam [KB14] with default parameters.