

# Layerwise and Dimensionwise Adaptive Local AMSMethod for Federated Learning

## Abstract

To be completed...

## 1 Introduction

A growing and important task while learning models on observed data, is the ability to train the latter over a large number of clients which could either be devices or distinct entities. In the paradigm of Federated Learning (FL) [3, 5], the focus of our paper, a central server orchestrates the optimization over those clients under the constraint that the data can neither be centralized nor shared among the clients. Most modern machine learning tasks can be casted as a large finite-sum optimization problem written as:

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta) \quad (1)$$

where  $n$  denotes the number of workers,  $f_i$  represents the average loss for worker  $i$  and  $\theta$  the global model parameter taking value in  $\Theta$  a subset of  $\mathbb{R}^d$ . While this formulation recalls that of distributed optimization, the core principle of FL is different that standard distributed paradigm.

FL currently suffers from two bottlenecks: communication efficiency and privacy. We focus on the former in this paper. While local updates, updates during which each client learn their local models, can reduce drastically the number of communication rounds between the central server and devices, new techniques must be employed to tackle this challenge. Some quantization [1, 6] or compression [4] methods allow to decrease the number of bits communicated at each round and are efficient method in a distributed setting. The other approach one can take is to accelerate the local training on each device and thus sending a better local model to the server at each round.

Under the important setting of heterogenous data, i.e. the data among each device can be distributed according to different distributions, current local optimization algorithms are perfectible. The most popular method for FL is using multiple local Stochastic Gradient Descent (SGD) steps in each device, sending those local models to the server that computes the average over those received local vector of parameters and broadcasts it back to the devices. This is called FEDAVG and has been introduced in [5].

In [2], the authors motivate the usage of adaptive gradient optimization methods as a better alternative to the standard SGD inner loop in FEDAVG. They propose an adaptive gradient method, namely LOCAL AMSGRAD, with communication cost sublinear in  $T$  that is guaranteed to converge to stationary points in  $\mathcal{O}(\sqrt{d/Tn})$ , where  $T$  is the number of iterations.

Based on recent progress in adaptive methods for accelerating the training procedure, see [7], we propose a variant of LOCAL AMSGRAD integrating dimensionwise and layerwise adaptive learning rate in each device's local update. Our contributions are as follows:

- We develop a novel optimization algorithm for federated learning, namely FED-LAMB, following a principled layerwise adaptation strategy to accelerate training of deep neural networks.
- theoretical results
- We exhibit the advantages of our method on several benchmarks supervised learning methods on both homogeneous and heterogeneous settings.

## 1.1 Related Work

Federated learning.

Adaptive gradient methods.

## 2 Layerwise and Dimensionwise Adaptive Methods

**Notations:** We denote by  $\theta$  the vector of parameters taking values in  $\mathbb{R}^d$ . For each layer  $\ell \in \llbracket L \rrbracket$ , where  $L$  is the total number of layers of the neural networks, and each coordinate  $j \in \llbracket p_\ell \rrbracket$  where  $p_\ell$  is the dimension per layer  $\ell$ , we note  $\theta^{\ell,j}$  its  $j$ th coordinate. The gradient of  $f$  with respect to  $\theta^\ell$  is denoted by  $\nabla_\ell f(\theta)$ . The index  $i \in \llbracket n \rrbracket$  denotes the index of the worker  $i$  in our federated framework.  $r$  and  $t$  are used as the round and local iteration numbers respectively. The smoothness per layer is denoted by  $L_\ell$  for each layer  $\ell \in \llbracket L \rrbracket$ .

### 2.1 Local AMS with LAMB

We propose a layerwise and dimensionwise local AMS algorithm in the following:

---

#### Algorithm 1 L&D LOCAL AMS FOR FEDERATED LEARNING

---

- 1: **Input:** parameter  $\beta_1, \beta_2$ , and learning rate  $\alpha_t$ .
  - 2: Init:  $\theta_0 \in \Theta \subseteq \mathbb{R}^d$ , as the global model shared by all devices and  $v_0 = \epsilon 1 \in \mathbb{R}^d$  and  $\bar{\theta}_0 = \frac{1}{n} \sum_{i=1}^n \theta_0$ .
  - 3: **for**  $r = 1$  to  $R$  **do**
  - 4:   Set  $\theta_{r,i}^0 = \bar{\theta}_{r-1}$
  - 5:   **parallel for device**  $d \in D^r$  **do:**
  - 6:    Compute stochastic gradient  $g_{r,i}$  at  $\theta_r$ .
  - 7:    **for**  $t = 1$  to  $T$  **do**
  - 8:       $m_{r,i}^t = \beta_1 m_{r-1,i}^{t-1} + (1 - \beta_1) g_{r,i}$ .
  - 9:       $\hat{m}_{r,i}^t = m_{r,i}^t / (1 - \beta_1^t)$ .
  - 10:      $v_r^{t,i} = \beta_2 v_{r-1,i}^t + (1 - \beta_2) g_{r,i}^2$ .
  - 11:      $\hat{v}_{r,i}^t = v_{r,i}^t / (1 - \beta_2^t)$ .
  - 12:      $\hat{v}_{r,i}^t = \max(\hat{v}_{r-1,i}^t, v_{r,i}^t)$ .
  - 13:     Compute ratio  $p_{r,i}^j = \frac{m_{r,i}^t}{\sqrt{\hat{v}_{r,i}^t + \epsilon}}$ .
  - 14:     Update local model for each layer  $\ell$ :
 
$$\theta_{r,i}^{\ell,t} = \theta_{r,i}^{\ell,t-1} - \alpha_r \phi(\|\theta_{r,i}^{\ell,t-1}\|)(p_{r,i}^j + \lambda \theta_{r,i}^{\ell,t-1}) / \|\theta_{r,i}^{\ell,t-1}\|$$
  - 15:   **end for**
  - 16:   Devices send local model  $\theta_{r,i}^T = [\theta_{r,i}^{\ell,T}]_{\ell=1}^L$  to the server
  - 17:   Server computes the averages of the local models  $\bar{\theta}_r^j = \frac{1}{n} \sum_{i=1}^n \theta_{r,i}^{j,T}$  and send it back to the devices.
  - 18: **end for**
- 

### 2.2 Finite time convergence bounds

In the context of nonconvex stochastic optimization for distributed devices, assume the following:

**H1.** For  $i \in \llbracket n \rrbracket$  and  $\ell \in \llbracket L \rrbracket$ ,  $f_i$  is  $L$ -smooth:  $\|\nabla f_i(\theta) - \nabla f_i(\vartheta)\| \leq L_\ell \|\theta^\ell - \vartheta^\ell\|$ .

We add some classical assumption in the unbiased stochastic optimization realm, on the gradient of the objective function:

**H2.** The stochastic gradient is unbiased for any iteration  $r > 0$ :  $\mathbb{E}[g_r] = \nabla f(\theta_r)$  and is bounded from above, i.e.,  $\|g_t\| \leq M$ .

**H3.** The variance of the stochastic gradient is bounded for any iteration  $r > 0$  and any dimension  $j \in \llbracket d \rrbracket$ :  $\mathbb{E}[|g_r^j - \nabla f(\theta_r)^j|^2] < \sigma^2$ .

**Case with  $T = 1$ ,  $\epsilon = 0$  and  $\lambda = 0$ :** Using H1, we have:

$$\begin{aligned} f(\bar{\vartheta}_{r+1}) &\leq f(\bar{\vartheta}_r) + \langle \nabla f(\bar{\vartheta}_r) | \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle + \sum_{\ell=1}^L \frac{L_\ell}{2} \|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2 \\ &\leq f(\bar{\vartheta}_r) + \sum_{\ell=1}^L \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j (\bar{\vartheta}_{r+1}^{\ell,j} - \bar{\vartheta}_r^{\ell,j}) + \sum_{\ell=1}^L \frac{L_\ell}{2} \|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2 \end{aligned} \quad (2)$$

Taking expectations on both sides leads to:

$$-\mathbb{E}[\langle \nabla f(\bar{\vartheta}_r) | \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle] \leq \mathbb{E}[f(\bar{\vartheta}_r) - f(\bar{\vartheta}_{r+1})] + \sum_{\ell=1}^L \frac{L_\ell}{2} \mathbb{E}[\|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2] \quad (3)$$

Yet, we observe that, using the classical intermediate quantity, used for proving convergence results of adaptive optimization methods, see [], we have:

$$\bar{\vartheta}_r = \bar{\theta}_r + \frac{\beta_1}{1 - \beta_1} (\bar{\theta}_r - \bar{\theta}_{r-1}) \quad (4)$$

where  $\bar{\theta}_r$  denotes the average of the local models at round  $r$ . Then for each layer  $\ell$ ,

$$\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell = \frac{1}{1 - \beta_1} (\bar{\theta}_{r+1}^\ell - \bar{\theta}_r^\ell) - \frac{\beta_1}{1 - \beta_1} (\bar{\theta}_r^\ell - \bar{\theta}_{r-1}^\ell) \quad (5)$$

$$= \frac{\alpha_r}{1 - \beta_1} \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\|p_{r,i}^\ell\|} p_{r,i}^\ell - \frac{\alpha_{r-1}}{1 - \beta_1} \sum_{i=1}^n \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\|p_{r-1,i}^\ell\|} p_{r-1,i}^\ell \quad (6)$$

$$= \frac{\alpha\beta_1}{1 - \beta_1} \sum_{i=1}^n \left( \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_{r,i}^t \|p_{r,i}^\ell\|}} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1,i}^t \|p_{r-1,i}^\ell\|}} \right) m_{r-1}^t + \alpha \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_{r,i}^t \|p_{r,i}^\ell\|}} g_{r,i} \quad (7)$$

where we have assumed a constant learning rate  $\alpha$ .

We note for all  $\theta \in \Theta$ , the majorant  $G > 0$  such that  $\phi(\|\theta\|) \leq G$ . Then, following (3), we obtain:

$$-\mathbb{E}[\langle \nabla f(\bar{\vartheta}_r) | \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle] \leq \mathbb{E}[f(\bar{\vartheta}_r) - f(\bar{\vartheta}_{r+1})] + \sum_{\ell=1}^L \frac{L_\ell}{2} \mathbb{E}[\|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2] \quad (8)$$

Developing the LHS of (8) using (5) leads to

$$\langle \nabla f(\bar{\vartheta}_r) | \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle = \sum_{\ell=1}^L \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j (\bar{\vartheta}_{r+1}^{\ell,j} - \bar{\vartheta}_r^{\ell,j}) \quad (9)$$

$$= \frac{\alpha\beta_1}{1 - \beta_1} \sum_{\ell=1}^L \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j \left[ \sum_{i=1}^n \left( \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_{r,i}^t \|p_{r,i}^\ell\|}} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1,i}^t \|p_{r-1,i}^\ell\|}} \right) m_{r-1}^t \right] \quad (10)$$

$$- \alpha \underbrace{\sum_{\ell=1}^L \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_{r,i}^t \|p_{r,i}^\ell\|}} g_{r,i}}_{=A_1} \quad (11)$$

Since we have that  $\|p_{r,i}^\ell\| \leq \sqrt{\frac{p_\ell}{1-\beta_2}}$  and  $\sqrt{v_{r,i}^t} \leq M$ , using H2, we develop the term  $A_1$  as follows:

$$A_1 \leq -\alpha \sum_{\ell=1}^L \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_{r,i}^t} \|p_{r,i}^\ell\|} g_{r,i} \quad (12)$$

$$\leq -\alpha \sum_{\ell=1}^L \sqrt{\frac{1-\beta_2}{M^2 p_\ell}} \sum_{i=1}^n \sum_{j=1}^{p_\ell} \phi(\|\theta_{r,i}^\ell\|) \nabla_\ell f(\bar{\vartheta}_r)^j g_{r,i}^{\ell,j} \quad (13)$$

$$- \alpha \sum_{\ell=1}^L \sum_{i=1}^n \sum_{j=1}^{p_\ell} \left( \phi(\|\theta_{r,i}^\ell\|) \nabla_\ell f(\bar{\vartheta}_r)^j \frac{p_{r,i}^\ell}{\|p_{r,i}^\ell\|} \right) 1 \left( \text{sign}(\nabla_\ell f(\bar{\vartheta}_r)^j) \neq \text{sign}(p_{r,i}^\ell) \right) \quad (14)$$

### 3 Numerical experiments

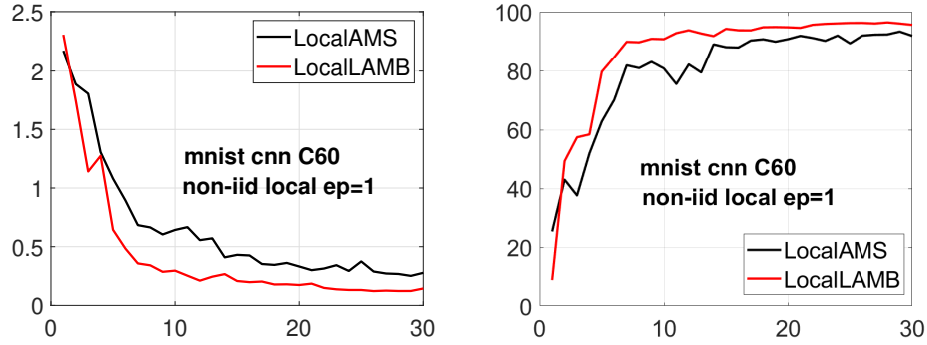


Figure 1: Test accuracy on CNN + MNIST. Non-iid data distribution.

### 4 Conclusion

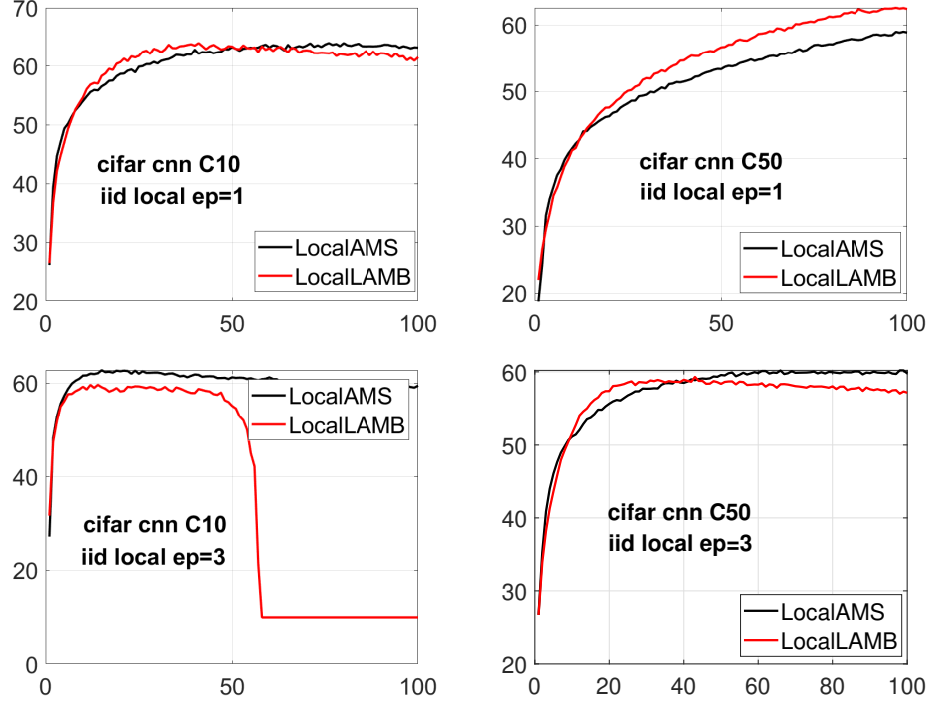


Figure 2: Test accuracy on CNN + CIFAR10. iid data distribution.

## References

- [1] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd: Communication-efficient sgd via gradient quantization and encoding. In *Advances in Neural Information Processing Systems*, pages 1709–1720, 2017.
- [2] Xiangyi Chen, Xiaoyun Li, and Ping Li. Toward communication efficient adaptive gradient method. In *ACM-IMS Foundations of Data Science Conference (FODS)*, Seattle, WA, 2020.
- [3] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- [4] Yujun Lin, Song Han, Huizi Mao, Yu Wang, and William J Dally. Deep gradient compression: Reducing the communication bandwidth for distributed training. *arXiv preprint arXiv:1712.01887*, 2017.
- [5] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pages 1273–1282. PMLR, 2017.
- [6] Jianqiao Wangni, Jiale Wang, Ji Liu, and Tong Zhang. Gradient sparsification for communication-efficient distributed optimization. In *Advances in Neural Information Processing Systems*, pages 1299–1309, 2018.
- [7] Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training bert in 76 minutes. *arXiv preprint arXiv:1904.00962*, 2019.

## A Appendix

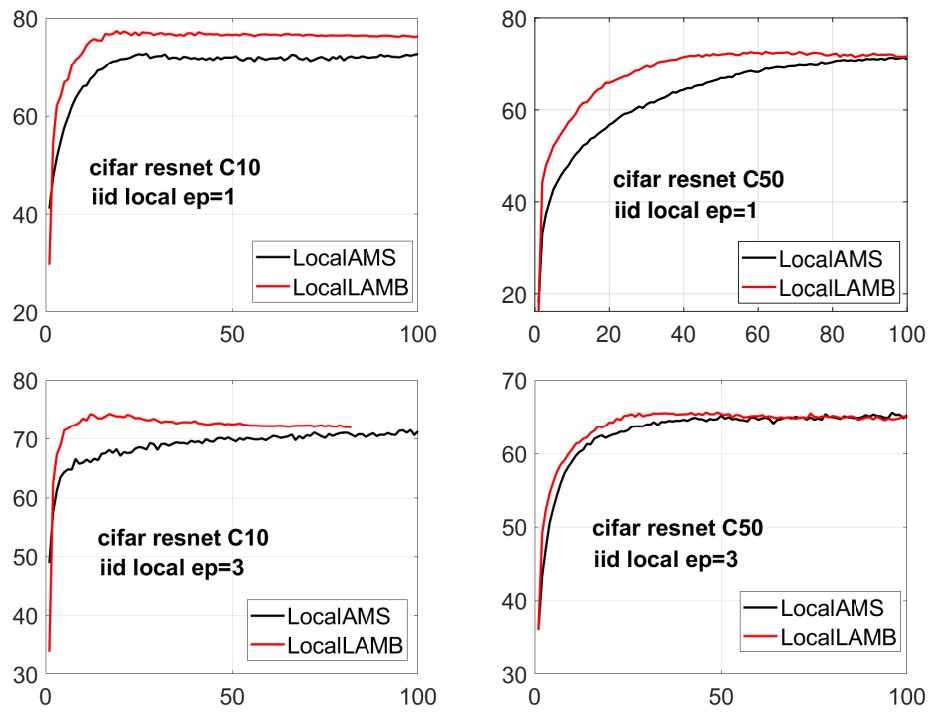


Figure 3: Test accuracy on ResNet + CIFAR10. iid data distribution.