

Proof for “Client-Edge-Cloud Hierarchical Federated Learning”

I. SYSTEM SETTINGS & ALGORITHMS

There are in total L edge servers indexed by ℓ with disjoint client sets $\{C^\ell\}_{\ell=1}^L$, N clients indexed by i with distributed datasets $\{\mathcal{D}_i^\ell\}_{i=1}^N$, \mathcal{D}^ℓ denotes the joint dataset under edge ℓ , $|\cdot|$ denotes the size of dataset. τ_1 is the number of local iterations before the clients upload its weights to the server, τ_2 is the number of aggregations performed on the edge before a global aggregation at the remote cloud server. Then with the HierFAVG algorithm in [1], the local model parameters $\mathbf{w}_i^\ell(t)$ evolves in the following way:

$$\mathbf{w}_i^\ell(t) = \begin{cases} \mathbf{w}_i^\ell(t-1) - \eta_t \nabla F_i^\ell(\mathbf{w}_i^\ell(t-1)) & t \mid \tau_1 \neq 0 \\ \frac{\sum_{i \in C^\ell} [\mathbf{w}_i^\ell(t-1) - \eta_t \nabla F_i^\ell(\mathbf{w}_i^\ell(t-1))]}{|\mathcal{D}^\ell|} & t \mid \tau_1 = 0 \\ \frac{\sum_{i=1}^N [\mathbf{w}_i^\ell(t-1) - \eta_t \nabla F_i^\ell(\mathbf{w}_i^\ell(t-1))]}{|\mathcal{D}|} & t \mid \tau_1 \tau_2 \neq 0 \\ \frac{\sum_{i=1}^N [\mathbf{w}_i^\ell(t-1) - \eta_t \nabla F_i^\ell(\mathbf{w}_i^\ell(t-1))]}{|\mathcal{D}|} & t \mid \tau_1 \tau_2 = 0 \end{cases} \quad (1)$$

II. DEFINITIONS & ASSUMPTIONS

The total T iterations of update is divided into K cloud intervals with length $\tau_1 \tau_2$ and $K \tau_2$ edge intervals with length τ_1 . We use $[p]$ to represent the iteration range $[(p-1)\tau_1, p\tau_1]$, $\{q\}$ to represent the iteration range $[(q-1)\tau_1 \tau_2, q\tau_1 \tau_2]$, so we have $\{q\} = \cup_p [p]$, $p = (q-1)\tau_2 + 1, (q-1)\tau_2 + 2, \dots, q\tau_2$.

To facilitate our analysis, we introduce several auxiliary sequences and virtual sequences:

- $F^\ell(\mathbf{w})$: Edge loss function

$$F^\ell(\mathbf{w}) = \frac{\sum_{i \in C^\ell} F_i(\mathbf{w})}{|\mathcal{D}^\ell|} \quad (2)$$

- $\bar{\mathbf{w}}^\ell(t)$: weighted average of all the $\mathbf{w}_i^\ell(t)$ under edge ℓ

$$\bar{\mathbf{w}}^\ell(t) = \frac{1}{|\mathcal{D}^\ell|} \sum_{i \in \mathcal{C}^\ell} |\mathcal{D}_i^\ell| \mathbf{w}_i^\ell(t) \quad (3)$$

- $\mathbf{w}(t)$: weighted average of all the $\mathbf{w}_i^\ell(t)$

$$\mathbf{w}(t) = \frac{1}{|\mathcal{D}|} \sum_{i=1}^N |\mathcal{D}_i^\ell| \mathbf{w}_i^\ell(t) \quad (4)$$

- $\mathbf{v}_{[p]}^\ell(t)$: Virtual edge centralized gradient descent sequence, defined on $[p]$ and evolves following (5), synchronized to $\bar{\mathbf{w}}^\ell(t)$ at the beginning of $[p]$

$$\begin{aligned} \mathbf{v}_{[p]}^\ell((p-1)\tau_1) &= \bar{\mathbf{w}}^\ell((p-1)\tau_1) \\ \mathbf{v}_{[p]}^\ell(t) &= \mathbf{v}_{[p]}^\ell(t-1) - \eta_t \nabla F^\ell(\mathbf{v}_{[p]}^\ell(t-1)) \end{aligned} \quad (5)$$

- $\tilde{\mathbf{v}}_{\{q\}}^\ell(t)$: Virtual cloud centralized gradient descent sequence, defined on $\{q\}$ and evolves following (6), synchronized to $\mathbf{w}(t)$ at the beginning of $\{q\}$

$$\begin{aligned} \tilde{\mathbf{v}}_{\{q\}}^\ell((q-1)\tau_1\tau_2) &= \mathbf{w}((q-1)\tau_1\tau_2) \\ \tilde{\mathbf{v}}_{\{q\}}^\ell(t) &= \tilde{\mathbf{v}}_{\{q\}}^\ell(t-1) - \eta_t \nabla F^\ell(\tilde{\mathbf{v}}_{\{q\}}^\ell(t-1)) \end{aligned} \quad (6)$$

- $\mathbf{u}_{\{q\}}(t)$: Virtual edge centralized gradient descent sequence, defined on $\{q\}$ and evolves following (7), synchronized to $\mathbf{w}(t)$ at the beginning of $\{q\}$

$$\begin{aligned} \mathbf{u}_{\{q\}}((q-1)\tau_1\tau_2) &= \mathbf{w}((q-1)\tau_1\tau_2) \\ \mathbf{u}_{\{q\}}(t) &= \mathbf{u}_{\{q\}}(t-1) - \eta_t \nabla F(\mathbf{u}_{\{q\}}(t-1)) \end{aligned} \quad (7)$$

We also have following assumptions on the loss functions.

Definition 1 (Gradient Divergence[2]) For any i and \mathbf{w} , define the gradient divergence between local loss function and edge loss function δ_i^ℓ as the upper bound of $\|\nabla F_i^\ell(\mathbf{w}) - \nabla F^\ell(\mathbf{w})\|$; the gradient divergence between edge loss function and global loss function Δ^ℓ as the upperbound of $\|\nabla F^\ell(\mathbf{w}) - \nabla F(\mathbf{w})\|$, i.e. ,

$$\begin{aligned} \|\nabla F_i^\ell(\mathbf{w}) - \nabla F^\ell(\mathbf{w})\| &\leq \delta_i^\ell \\ \|\nabla F^\ell(\mathbf{w}) - \nabla F(\mathbf{w})\| &\leq \Delta^\ell \end{aligned}$$

Define $\delta = \frac{\sum_{i=1}^N |\mathcal{D}_i^\ell| \delta_i^\ell}{|\mathcal{D}|}$, $\Delta = \frac{\sum_{\ell=1}^L |\mathcal{D}^\ell| \Delta^\ell}{|\mathcal{D}|} = \frac{\sum_{i=1}^N |\mathcal{D}_i^\ell| \Delta^\ell}{|\mathcal{D}|}$

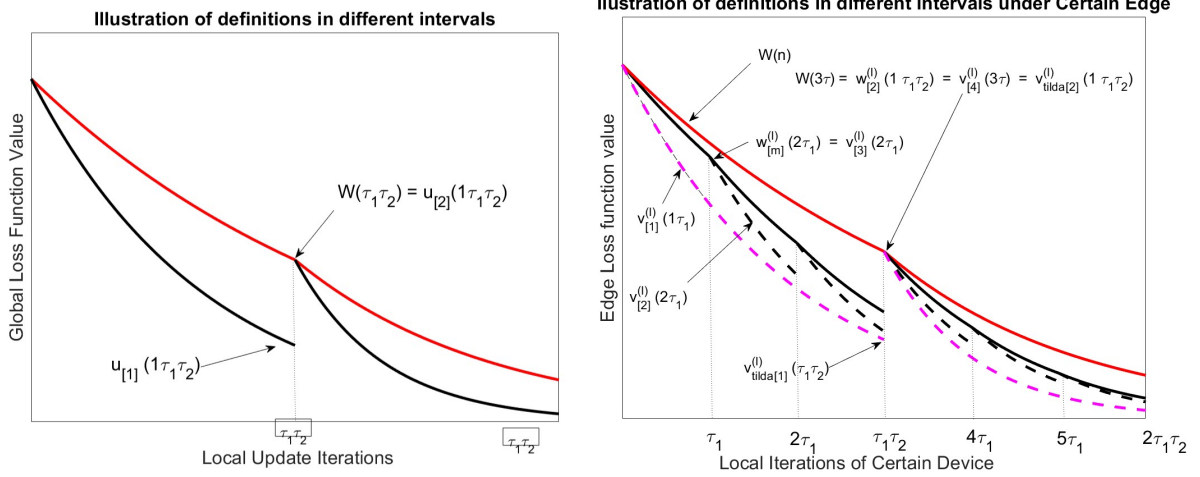


Fig. 1: Illustrations of the virtual sequences and auxiliary variables($\tau_2 = 3$)

Assumption 1 (Lipschitz-continuous) For any i , $F_i^\ell(\mathbf{w})$ is ρ -continuous, i.e., $\|F_i^\ell(\mathbf{w}) - F_i^\ell(\mathbf{w}')\| \leq \rho\|\mathbf{w} - \mathbf{w}'\|$

Assumption 2 (Lipschitz-smooth) For any i , $F_i^\ell(\mathbf{w})$ is β -smooth, i.e., $\|\nabla F_i^\ell(\mathbf{w}) - \nabla F_i^\ell(\mathbf{w}')\| \leq \beta\|\mathbf{w} - \mathbf{w}'\|$

Assumption 3 For any i , $F_i^\ell(\mathbf{w})$ is convex.

It can be easily inferred that $F^\ell(\mathbf{w})$ and $F(\mathbf{w})$ are all i) ρ -Lipschitz under Assumption 1 ii) β -smooth under Assumption 2 iii) convex under Assumption 3

III. PROOF OF LEMMA 2 & LEMMA 3 IN [1]

Lemma 2 (Convex) With Assumptions 1, 2, and 3, for any cloud aggregation interval $\{q\}$ with a fixed step size η_q , and $t \in \{q\}$, we have

$$\|\mathbf{w}(t) - \mathbf{u}_{\{q\}}(t)\| \leq G(t, \eta_q) \quad (8)$$

where

$$\begin{aligned} G(t, \eta_q) = & h(t - (q-1)\tau_1\tau_2, \Delta, \eta_q) + h(t - ((q-1)\tau_2 + p(t) - 1)\tau_1, \delta, \eta_q) \\ & + \frac{\tau_1}{2}(p^2(t) + p(t) - 2)h(\tau_1, \delta, \eta_q) \end{aligned} \quad (9)$$

$$\begin{aligned} h(x, \delta, \eta) = & \frac{\delta}{\beta}((\eta\beta + 1)^x - 1) - \eta\beta x \\ p(x) = & \lceil \frac{x}{\tau_1} - (q-1)\tau_2 \rceil \end{aligned} \quad (10)$$

Proof. From Eqn. (4) and the HierFAVG training algorithm Eqn. (1), we have

$$\mathbf{w}(t) = \frac{1}{|\mathcal{D}|} \sum_{i=1}^N \mathbf{w}_i^\ell(t) = \mathbf{w}(t-1) - \frac{\eta_q}{|\mathcal{D}|} \sum_{i=1}^N \nabla F_i^\ell(\mathbf{w}_i^\ell(t-1)) \quad (11)$$

Thus,

$$\begin{aligned} \|\mathbf{w}(t) - \mathbf{u}_{\{q\}}(t)\| - \|\mathbf{w}(t-1) - \mathbf{u}_{\{q\}}(t-1)\| &\leq \frac{\eta_q}{|\mathcal{D}|} \sum_{i=1}^N |D_j^\ell| \times \|\nabla F_i^\ell(\mathbf{w}_j^\ell(t-1)) - \nabla F_i^\ell(\mathbf{u}_{\{q\}}(t-1))\| \\ &\leq \frac{\eta_q \beta}{|\mathcal{D}|} \sum_{i=1}^N |D_j^\ell| \times \|\mathbf{w}_j^\ell(t-1) - \mathbf{u}_{\{q\}}(t-1)\|. \end{aligned} \quad (12)$$

Since we have $\mathbf{u}_{\{q\}}((q-1)\tau_1\tau_2) = \mathbf{w}((q-1)\tau_1\tau_2)$, thus

$$\begin{aligned} \|\mathbf{w}(t) - \mathbf{u}_{\{q\}}(t)\| &= \sum_{y=(q-1)\tau_1\tau_2+1}^t \|\mathbf{w}(y) - \mathbf{u}_{\{q\}}(y)\| - \|\mathbf{w}(t-1) - \mathbf{u}_{\{q\}}(t-1)\| \\ &\leq \frac{\eta_q \beta}{|\mathcal{D}|} \sum_{y=(q-1)\tau_1\tau_2+1}^t \sum_{i=1}^N |D_j^\ell| \times \|\mathbf{w}_j^\ell(y-1) - \mathbf{u}_{\{q\}}(y-1)\|. \end{aligned} \quad (13)$$

Next, we will show how to bound $\|\mathbf{w}_j^\ell(t) - \mathbf{u}_{\{q\}}(t)\|$. Using triangle inequality, we have:

$$\|\mathbf{w}_j^\ell(t) - \mathbf{u}_{\{q\}}(t)\| \leq \|\mathbf{w}_j^\ell(t) - \mathbf{v}_{[p]}^\ell(t)\| + \|\mathbf{v}_{[p]}^\ell(t) - \tilde{\mathbf{v}}_{\{q\}}^\ell(t)\| + \|\tilde{\mathbf{v}}_{\{q\}}^\ell(t) - \mathbf{u}_{\{q\}}(t)\| \quad (14)$$

Using Lemma 3 in [2], we can get:

$$\begin{aligned} \|\mathbf{w}_j^\ell(t) - \mathbf{v}_{[p]}^\ell(t)\| &\leq g_i^\ell(t - (p-1)\tau_1, \delta_i^\ell, \eta_q) \\ \|\tilde{\mathbf{v}}_{\{q\}}^\ell(t) - \mathbf{u}_{\{q\}}(t)\| &\leq g^\ell(t - (q-1)\tau_1\tau_2, \Delta^\ell, \eta_q) \end{aligned} \quad (15)$$

where function $g(x; \delta)$ is defined as $g_i^\ell(x; \delta) = \frac{\delta_i^\ell}{\beta}((1 + \eta_q \beta)^x - 1)$, function $g^\ell(x; \Delta)$ is defined as $g^\ell(x; \Delta) = \frac{\Delta^\ell}{\beta}((1 + \eta_q \beta)^x - 1)$.

By summing it over y and i in Eqn. (13), we have

$$\begin{aligned} \frac{\eta_q \beta}{|\mathcal{D}|} \sum_{y=(q-1)\tau_1\tau_2+1}^t \sum_{i=1}^N |D_j^\ell| \times \|\mathbf{w}_j^\ell(y) - \mathbf{v}_{[p]}^\ell(y)\| &\leq \frac{\eta_q \beta}{|\mathcal{D}|} \sum_{y=(q-1)\tau_1\tau_2+1}^t \sum_{i=1}^N |D_j^\ell| \times g_i^\ell(y - (p-1)\tau_1, \delta_i^\ell, \eta_q) \\ &\leq h(t - ((q-1)\tau_2 + p(t) - 2)\tau_1, \delta, \eta_q) \end{aligned} \quad (16)$$

$$\begin{aligned} \frac{\eta_q \beta}{|\mathcal{D}|} \sum_{y=(q-1)\tau_1\tau_2+1}^t \sum_{i=1}^N |D_j^\ell| \times \|\tilde{\mathbf{v}}_{\{q\}}^\ell(y) - \mathbf{u}_{\{q\}}(y)\| &\leq \frac{\eta_q \beta}{|\mathcal{D}|} \sum_{y=(q-1)\tau_1\tau_2+1}^t \sum_{i=1}^N |D_j^\ell| \times g^\ell(y - (q-1)\tau_1\tau_2, \Delta^\ell, \eta_q) \\ &\leq h(t - (q-1)\tau_1\tau_2, \Delta, \eta_q) \end{aligned} \quad (17)$$

where $p(t) = \lceil \frac{t}{\tau_1} - (q-1)\tau_2 \rceil$, $p(t)$ is defined the index of edge aggregation in the cloud interval $\{q\}$.

For $\|v_{[p]}^\ell(t) - \tilde{v}_{\{q\}}^\ell(t)\|$, we will show that it is bounded by a step function and increases at $t = (q-1)\tau_1\tau_2 + p\tau_1, p = 1, \dots, \tau_2 - 1$, which means that for $t \in [p]$,

$$\|v_{[p]}^\ell(t) - \tilde{v}_{\{q\}}^\ell(t)\| \leq \|v_{[p]}^\ell((p-1)\tau_1) - \tilde{v}_{\{q\}}^\ell((p-1)\tau_1)\| \quad (18)$$

From Eqn. (5) and (6), we can see that when $t \in [p]$ $v_{[p]}^\ell(t)$ and $\tilde{v}_{\{q\}}^\ell(t)$, the evolution of these two parameters are equal to performing gradient descent on the same loss function from different initialization point. A convex and smooth function $f(x)$ satisfies the following:

$$\begin{aligned} \langle \nabla f(x) - \nabla f(y), x - y \rangle &\geq \mu \|x - y\|_2^2, \\ \langle \nabla f(x) - \nabla f(y), x - y \rangle &\geq \frac{1}{\beta} \|\nabla f(x) - \nabla f(y)\|_2^2 \end{aligned} \quad (19)$$

where $\mu \geq 0$. Thus, it satisfies the regularity condition, i.e.,

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \frac{1}{2} \mu \|x - y\|_2^2 + \frac{1}{2\beta} \|\nabla f(x) - \nabla f(y)\|_2^2 \quad (20)$$

Thus for $v_{[p]}^\ell(t+1)$ and $\tilde{v}_{\{q\}}^\ell(t+1)$,

$$\begin{aligned} &\left\| v_{[p]}^\ell(t+1) - \tilde{v}_{\{q\}}^\ell(t+1) \right\|_2^2 - \left\| v_{[p]}^\ell(t) - \tilde{v}_{\{q\}}^\ell(t) \right\|_2^2 \\ &= \left\| \nabla F^\ell(v_{[p]}^\ell(t)) - \nabla F^\ell(\tilde{v}_{\{q\}}^\ell(t)) \right\|_2^2 - 2\eta \langle v_{[p]}^\ell(t) - \tilde{v}_{\{q\}}^\ell(t), \nabla F^\ell(v_{[p]}^\ell(t)) - \nabla F^\ell(\tilde{v}_{\{q\}}^\ell(t)) \rangle \\ &\leq 0 \end{aligned} \quad (21)$$

when $\eta \leq \frac{1}{\beta}$. Thus, Eqn. (18) is proved. The incremental value of the step function is $h(\tau_1, \delta, \eta_q)$.

Thus,

$$\|v_{[p]}^\ell(t) - \tilde{v}_{\{q\}}^\ell(t)\| \leq p(t) * h(\tau_1, \delta, \eta_q) \quad (22)$$

summing Eqn. 22, we have

$$\frac{\eta_q \beta}{|\mathcal{D}|} \sum_{y=(q-1)\tau_1\tau_2+1}^t \sum_{i=1}^N |D_j^\ell| \times \|v_{[p]}^\ell(t) - \tilde{v}_{\{q\}}^\ell(t)\| \leq \frac{\tau_1}{2} (p^2(t) + p(t) - 2) h(\tau_1, \delta, \eta_q) \quad (23)$$

By adding the bounded value of the three terms in Eqn. (14), the bound in Lemma 2 is derived. \square

Lemma 3 (Non-convex) For any i , assuming $f_i(w)$ is β -smooth, for any cloud interval $\{q\}$ with

step size η_q , we have

$$\|\mathbf{w}(t) - \mathbf{u}_{\{q\}}(t)\| \leq G_{nc}(\kappa_1 \kappa_2, \eta_q)$$

where

$$G_{nc}(\kappa_1 \kappa_2, \eta_q) = h(\kappa_1 \kappa_2, \Delta, \eta_q) + \kappa_1 \kappa_2 \frac{(1 + \eta_q \beta)^{\kappa_1 \kappa_2} - 1}{(1 + \eta_q \beta)^{\kappa_1} - 1} h(\kappa_1, \delta, \eta_q) + h(\kappa_1, \delta, \eta_q),$$

$$h(x, \delta, \eta) = \frac{\delta}{\beta} ((\eta \beta + 1)^x - 1) - \eta \beta x.$$

Proof. The proof process is the same as the convex case except for bounding the $\|v_{[p]}^\ell(t) - \tilde{v}_{\{q\}}^\ell(t)\|$, since the regularity condition is not satisfied for non-convex functions. For non-convex loss functions, we have

$$\begin{aligned} \|v_{[p]}^\ell(t+1) - \tilde{v}_{\{q\}}^\ell(t+1)\| &= \|v_{[p]}^\ell(t) - \tilde{v}_{\{q\}}^\ell(t) - (\eta_q \nabla F^\ell(v_{[p]}^\ell(t+1)) - \eta_q \nabla F^\ell \tilde{v}_{\{q\}}^\ell(t+1))\| \\ &\leq \|v_{[p]}^\ell(t) - \tilde{v}_{\{q\}}^\ell(t)\| + \eta_q \beta \|v_{[p]}^\ell(t) - \tilde{v}_{\{q\}}^\ell(t)\| \\ &\leq (1 + \eta_q \beta) \|v_{[p]}^\ell(t) - \tilde{v}_{\{q\}}^\ell(t)\| \\ &\leq (1 + \eta_q \beta)^{t-(p-1)\tau_1} \|v_{[p]}^\ell((p-1)\tau_1) - \tilde{v}_{\{q\}}^\ell((p-1)\tau_1)\| \end{aligned} \tag{24}$$

□

REFERENCES

- [1] L. Liu, J. Zhang, S. H. Song, and K. B. Letaief, "Client-Edge-Cloud Hierarchical Federated Learning," in *IEEE International Conference on Communications (ICC)*, 2020, pp. 1–6.
- [2] S. Wang, T. Tuor, T. Salonidis, K. K. Leung, C. Makaya, T. He, and K. Chan, "When edge meets learning: Adaptive control for resource-constrained distributed machine learning," in *IEEE INFOCOM 2018-IEEE Conference on Computer Communications*. IEEE, 2018, pp. 63–71.