

Enabling Communication-efficient and Robust Federated Learning over Packet Lossy Networks via Random Interleaved Vector Quantization

IEEE International Conference on Multimedia and Expo (ICME) 2025

APPENDIX

Proof of Theorem 1. Define $\mathbf{g} = (g_1, g_2, \dots, g_n)$ as raw gradient, the received lossy gradient compressed by RIVQ as $\tilde{\mathbf{g}}$, erasure rate as r and the compression error as $\mathbf{e} = (e_1, e_2, \dots, e_n)$, then the cosine distance between \mathbf{g} and $\tilde{\mathbf{g}}$ satisfies:

$$\begin{aligned} \tilde{d}_{\cos}(\tilde{\mathbf{g}}, \mathbf{g}) &\stackrel{(a)}{=} 1 - \frac{\langle \mathbf{g}, \tilde{\mathbf{g}} \rangle}{\|\mathbf{g}\| \|\tilde{\mathbf{g}}\|} \\ &\stackrel{(b)}{=} 1 - \frac{\|\tilde{\mathbf{g}}\|^2}{\|\mathbf{g}\| \sqrt{\sum_{i \in S_e^c} (g_i^2 + e_i^2 + 2g_i e_i)}} \\ &\stackrel{(c)}{\leq} 1 - \frac{(1-r)\|\mathbf{g}\|^2}{\|\mathbf{g}\| \sqrt{\sum_{i \in S_e^c} 2(g_i^2 + e_i^2)}} \\ &\stackrel{(d)}{\leq} 1 - \frac{(1-r)\|\mathbf{g}\|^2}{\|\mathbf{g}\| \sqrt{2(1-r) \sum_{i=1}^n (g_i^2 + e_i^2)}} \\ &\stackrel{(e)}{\leq} 1 - \frac{\sqrt{(1-r)\|\mathbf{g}\|^2}}{\sqrt{2(\|\mathbf{e}\|^2 + \|\mathbf{g}\|^2)}} \end{aligned} \quad (1)$$

where inequality (b) and (c) follow from the fact that $\langle \mathbf{g}, \tilde{\mathbf{g}} \rangle = \sum_{i \in S_e^c} g_i^2 = \|\tilde{\mathbf{g}}\|^2$, $\|\tilde{\mathbf{g}}\|^2 = \sum_{i \in S_e^c} g_i^2 \approx (1-r) \sum_{i=1}^n g_i^2$ and Young's inequality. Inequality (d) holds under $\sum_{i \in S_e^c} e_i^2 \approx (1-r) \sum_{i=1}^n e_i^2 = (1-r)\|\mathbf{e}\|^2$. S_e^c refers to the set of received gradient entries, i.e., complement of erasure entry set S_e . \square

Proof of Theorem 2. Supposing global loss function $F(\mathbf{w}_t) = \frac{1}{K} \sum_{k=1}^K f_k(\mathbf{w}_t)$ satisfies L -smooth (Assumption 1), we have

$$\begin{aligned} &F(\mathbf{w}_{t+1}) - F(\mathbf{w}_t) \\ &\stackrel{(a)}{\leq} \langle \nabla F(\mathbf{w}_t), \mathbf{w}_{t+1} - \mathbf{w}_t \rangle + \frac{L}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 \\ &\stackrel{(b)}{\leq} -\eta \langle \nabla F(\mathbf{w}_t), \tilde{\mathbf{g}}_t \rangle + \frac{L\eta^2}{2} \|\tilde{\mathbf{g}}_t\|^2 \\ &\stackrel{(c)}{\leq} -\eta \langle \nabla F(\mathbf{w}_t), (\tilde{\mathbf{g}}_t - \mathbf{g}_t) + \mathbf{g}_t \rangle + \frac{L\eta^2}{2} \|(\tilde{\mathbf{g}}_t - \mathbf{g}_t) + \mathbf{g}_t\|^2 \\ &\stackrel{(d)}{\leq} -\eta \|\nabla F(\mathbf{w}_t)\|^2 + \eta \|\nabla F(\mathbf{w}_t)\| \|\tilde{\mathbf{g}}_t - \mathbf{g}_t\| \\ &\quad + L\eta^2 \|(\tilde{\mathbf{g}}_t - \mathbf{g}_t)\|^2 + L\eta^2 \|\mathbf{g}_t\|^2 \\ &\stackrel{(e)}{\leq} -\eta \|\nabla F(\mathbf{w}_t)\|^2 + \eta \|\nabla F(\mathbf{w}_t)\|^2 \sqrt{2\tilde{d}_{\cos}(\tilde{\mathbf{g}}, \mathbf{g})} \\ &\quad + 2L\eta^2 \|\nabla F(\mathbf{w}_t)\|^2 \tilde{d}_{\cos}(\tilde{\mathbf{g}}, \mathbf{g}) + L\eta^2 \|\nabla F(\mathbf{w}_t)\|^2 \\ &\stackrel{(f)}{\leq} \|\nabla F(\mathbf{w}_t)\|^2 [\eta^2 L - \eta + \eta \sqrt{2\tilde{d}_{\cos}(\mathbf{g}, \tilde{\mathbf{g}})} + 2\eta^2 L \tilde{d}_{\cos}(\tilde{\mathbf{g}}, \mathbf{g})] \end{aligned} \quad (2)$$

where inequality (a) holds under Taylor's theorem [1] and L -smooth (Assumption 1); (b) commences from $\mathbf{w}_{t+1} = \mathbf{w}_t -$

$\eta \tilde{\mathbf{g}}_t$ and $\tilde{\mathbf{g}}_t = \mathbf{g}_t + \mathbf{e}_t$; (d) is derived by Cauchy-Schwarz inequality and triangle inequality; (e) originates from the fact that $\|\mathbf{g} - \tilde{\mathbf{g}}\| \approx \|\mathbf{g}_t\| \sqrt{2\tilde{d}_{\cos}(\tilde{\mathbf{g}}_t, \mathbf{g}_t)}$.

Then, we compute the telescoping sums of inequality (2) over $t = 1, 2, \dots, T$ and obtain:

$$\begin{aligned} F(\mathbf{w}^*) - F(\mathbf{w}_1) &\leq F(\mathbf{w}_{T+1}) - F(\mathbf{w}_1) \leq [\eta - \eta^2 L \\ &\quad - \eta \sqrt{2\tilde{d}_{\cos}(\tilde{\mathbf{g}}, \mathbf{g})} - 2\eta^2 L \tilde{d}_{\cos}(\tilde{\mathbf{g}}, \mathbf{g})] \sum_{t=0}^{T-1} \|\nabla F(\mathbf{w}_t)\|^2 \end{aligned} \quad (3)$$

By substituting $\eta = \sqrt{\frac{K}{LT}}$ into inequality (3), then multiplying by $\frac{1}{T}$, taking expectations on its both sides and rewriting, we get desired results (omit some subscripts for brevity). \square

Proof of the Effectiveness of Gradient Random Interleaving. Define the raw gradient as $\mathbf{g} = (g_1, g_2, \dots, g_n)$, \mathbf{g} is equally partitioned into M packets, i.e., $\mathbf{g} = (P_1, P_2, \dots, P_M)$, P_i/g_i is described as P_i / \mathcal{P}_i when packet i is erased.

(1). *Gradient Entry Erasure after Random Interleaving*

For g_i , whether it is erased or not depends on which packet is assigned to after interleaving. Designate the packet index of g_i after interleaving as J_i , where:

$$P(J_i = j) = \frac{1}{M}, \forall j \in 1, 2, \dots, M$$

g_i is erased only when the packet $J_i = j$ is erased, hence:

$$P(\mathcal{P}_i) = \sum_{i=1}^M \frac{1}{M} \cdot P(\mathcal{P}_j) = r$$

(2). *Independence of Gradient Entry Erasure Events*

For two arbitrary gradient entries, random interleaving enables that whether or not they are assigned to the same packet is independent. Therefore, they are conditionally independent as well when it comes to packet erasure.

$$P(g_i \in P_j | g_{i'} \in P_j) = P(g_i \in P_j)$$

$$P(\mathcal{P}_i | \mathcal{P}_{i'}) = P(\mathcal{P}_i)$$

$$P(g_i | g_{i'}) = P(g_i)$$

In summary, through eliminating the dependencies among gradient entries, random interleaving renders packet erasure equivalent to random sparsification successfully. \square

Assumption 1 (L -Smooth [2]). *The loss function f in clients' local updating is Lipschitz-smooth with a constant L and uniformly bounded by $f(\mathbf{w}^*)$, that is, $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,*

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$$

REFERENCES

- [1] Gilbert Strang, *Calculus*, vol. 1, SIAM, 1991.
- [2] Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang, “On the convergence of FedAvg on Non-IID data,” in *ICLR*, 2020.