

A APPENDIX

A.1 Proof of Lemma 1.

In order to demonstrate the convergence of our update scheme, we make the following assumptions:

- Lipschitzian gradient: $f(\mathbf{x})$ is L -Lipschitz smooth, i.e., $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$, $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.
- Bounded variance: the variance of the stochastic gradient, denoted by ξ , is bounded such that $\mathbb{E}[\|\xi\|^2] \leq \sigma^2$.

Due to the L -Lipschitz condition satisfied by $f(\mathbf{x})$, the derivative of the formula $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$ with respect to \mathbf{x} yields $\nabla^2 f(\mathbf{x}) \leq L$. Therefore, we can obtain the second-order Taylor expansion of $f(\mathbf{x})$ around $\nabla^2 f(\mathbf{x}) \leq L$.

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2 \quad (1)$$

Therefore, we can deduce the following proof.

$$\begin{aligned} & \mathbb{E}f(\mathbf{x}_{t+1}) - \mathbb{E}f(\mathbf{x}_t) \\ & \leq \mathbb{E} \langle \mathbf{x}_{t+1} - \mathbf{x}_t, \nabla f(\mathbf{x}_t) \rangle + \frac{L}{2} \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ & = \mathbb{E} \langle -\gamma \nabla g(\mathbf{x}_t), \nabla g(\mathbf{x}_t) + p\xi_t \rangle + \frac{LY^2}{2} \mathbb{E} \|\nabla g(\mathbf{x}_t) + p\xi_t\|^2 \quad (\nabla g(\mathbf{x}) = \nabla f(\mathbf{x}) - \xi \text{ and } \mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \nabla g(\mathbf{x}_t)) \\ & \leq -\gamma \mathbb{E} \langle \nabla g(\mathbf{x}_t), \nabla g(\mathbf{x}_t) + p\xi_t \rangle + \frac{LY^2}{2} \mathbb{E} \|\nabla g(\mathbf{x}_t) + p\xi_t\|^2 \\ & \leq -\gamma \mathbb{E} \|\nabla g(\mathbf{x}_t)\|^2 + \frac{LY^2}{2} p^2 \mathbb{E} \|\nabla g(\mathbf{x}_t)\|^2 + \frac{LY^2}{2} p^2 \mathbb{E} \xi_t^2 \quad (\mathbb{E}[\xi] = 0) \\ & \leq (-\gamma + \frac{LY^2}{2}) \mathbb{E} \|\nabla g(\mathbf{x}_t)\|^2 + \frac{LY^2}{2} p^2 \sigma^2 \quad (\mathbb{E}[\|\xi\|^2] \leq \sigma^2) \end{aligned}$$

Using the method of telescoping sum, we can obtain the following equations.

$$\begin{aligned} \mathbb{E}f(\mathbf{x}_t) - \mathbb{E}f(\mathbf{x}_0) & \leq \left(\frac{LY^2}{2} - \gamma\right) \sum_{t=0}^{t-1} \mathbb{E} \|\nabla g(\mathbf{x}_t)\|^2 + \frac{LY^2}{2} p^2 \sigma^2 T \\ \left(\gamma - \frac{LY^2}{2}\right) \sum_{t=0}^{t-1} \mathbb{E} \|\nabla g(\mathbf{x}_t)\|^2 & \leq \mathbb{E}f(\mathbf{x}_0) - \mathbb{E}f(\mathbf{x}_t) + \frac{LY^2}{2} p^2 \sigma^2 T \\ \left(\gamma - \frac{LY^2}{2}\right) \frac{1}{T} \sum_{t=0}^{t-1} \mathbb{E} \|\nabla g(\mathbf{x}_t)\|^2 & \leq \frac{\mathbb{E}f(\mathbf{x}_0) - \mathbb{E}f(\mathbf{x}_t)}{T} + \frac{LY^2}{2} p^2 \sigma^2 \\ \frac{1}{T} \sum_{t=0}^{t-1} \mathbb{E} \|\nabla g(\mathbf{x}_t)\|^2 & \leq \frac{2L(\mathbb{E}f(\mathbf{x}_0) - \mathbb{E}f(\mathbf{x}_t))}{T} + p^2 \sigma^2 \end{aligned}$$

Lemma 1 suggests that by setting the probability of selecting hybrid batches for training to $p = 1/\sqrt{T}$, we achieve a convergence rate of $O(1/T)$, equivalent to that of full-vertex batch training.

$$\frac{1}{T} \sum_{t=0}^{t-1} \mathbb{E} \|\nabla g(\mathbf{x}_t)\|^2 \leq \frac{2L(\mathbb{E}f(\mathbf{x}_0) - \mathbb{E}f(\mathbf{x}_t)) + \sigma^2}{T} \quad (2)$$

Setting $p = 1/\sqrt[4]{T}$, however, leads to a convergence rate of $O(1/\sqrt{T})$, akin to mini-batch training.

$$\frac{1}{T} \sum_{t=0}^{t-1} \mathbb{E} \|\nabla g(\mathbf{x}_t)\|^2 \leq \frac{2L(\mathbb{E}f(\mathbf{x}_0) - \mathbb{E}f(\mathbf{x}_t))}{T} + \frac{\sigma^2}{\sqrt{T}} \quad (3)$$