

New Non-Convex Analysis of Local SGD and SCAFFOLD

Ruichen Luo¹

Joint work with

Sebastian Stich²

Samuel Horvath³

Martin Takac³

¹Institute of Science and Technology Austria (ISTA)

²CISPA Helmholtz Center for Information Security

³Mohamed bin Zayed University of Artificial Intelligence (MBZUAI)

MLO Group Meeting, St. Intberg, Germany, Oct 23, 2024

Distributed non-convex optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$$

- ▶ $\mathbf{x} \in \mathbb{R}^d$ is the optimization variable
- ▶ f is the global objective function, bounded below by f^*
- ▶ n is the number of workers
- ▶ f_i is the local objective function distributed to the i th worker, and f_i has L -Lipschitz continuous gradient, for each $i \in \{1, \dots, n\}$
- ▶ The local objective functions are *heterogeneous* in general, i.e., $f_i \neq f$

Stochastic oracle and intermittent communication

T subsequent queries to a *fully stochastic oracle* \mathcal{SO} :

- ▶ The workers input $(\mathbf{x}_t^1, \dots, \mathbf{x}_t^n) \in \mathbb{R}^{d \times n}$
- ▶ The \mathcal{SO} outputs $(G_1(\mathbf{x}_t^1, \xi_t^1), \dots, G_n(\mathbf{x}_t^n, \xi_t^n)) \in \mathbb{R}^{d \times n}$
- ▶ $\{\xi_t^i : 0 \leq t \leq T - 1\}$ are i.i.d. random variables
- ▶ Assume $\mathbb{E}_{\xi_t^i}[G_i(\mathbf{x}, \xi_t^i)] = \nabla f_i(\mathbf{x})$, $\mathbb{E}_{\xi_t^i} \|G_i(\mathbf{x}, \xi_t^i) - \nabla f_i(\mathbf{x})\|_2^2 \leq \sigma^2$

The worker communicates after every τ iterations:

- ▶ Assume T is a multiple of τ

Notations:

- ▶ $(\mathbf{g}_t^1, \dots, \mathbf{g}_t^n) = (G_1(\mathbf{x}_t^1, \xi_t^1), \dots, G_n(\mathbf{x}_t^n, \xi_t^n))$
- ▶ $\bar{\mathbf{x}}_t = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_t^i$
- ▶ $f^* = \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$
- ▶ $\Delta = f(\bar{\mathbf{x}}_0) - f^*$

Minibatch SGD vs. Local SGD/SCAFFOLD

Initialization: $\mathbf{x}_0^1 = \dots = \mathbf{x}_0^n$

MbSGD.

$$\mathbf{x}_{t+1}^i = \begin{cases} \bar{\mathbf{x}}_{t-\tau+1} - \frac{\eta}{n} \sum_{j=1}^n \sum_{k=0}^{\tau-1} \mathbf{g}_{t-k}^j, & \text{if } t+1 \text{ is a multiple of } \tau, \\ \mathbf{x}_t^i, & \text{otherwise} \end{cases}$$

LocalSGD.

$$\mathbf{x}_{t+1}^i = \begin{cases} \bar{\mathbf{x}}_{t-\tau+1} - \frac{\eta}{n} \sum_{j=1}^n \sum_{k=0}^{\tau-1} \mathbf{g}_{t-k}^j, & \text{if } t+1 \text{ is a multiple of } \tau \\ \mathbf{x}_t^i - \eta \mathbf{g}_t^i, & \text{otherwise} \end{cases}$$

SCAFFOLD.

(next page)

Minibatch SGD vs. Local SGD/SCAFFOLD

Algorithm 1 SCAFFOLD

```
1: for  $r = 0, 1, \dots, R - 1$  do
2:   for  $i \in [n]$  do in parallel
3:     for  $k = 0, 1, \dots, \tau - 1$  do
4:        $\mathbf{x}_{2r\tau+k+1}^i = \mathbf{x}_{2r\tau+k}^i$ 
5:     end for
6:      $\hat{\mathbf{g}}_{(r\tau)}^i = \frac{1}{\tau} \sum_{k=0}^{\tau-1} \mathbf{g}_{2r\tau+k}^i$ 
7:   end for
8:   Compute and broadcast:  $\hat{\mathbf{g}}_{(r\tau)} = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{g}}_{(r\tau)}^i$ 
9:   for  $i \in [n]$  do in parallel
10:    for  $k = \tau, \tau + 1, \dots, 2\tau - 2$  do
11:       $\mathbf{x}_{2r\tau+k+1}^i = \mathbf{x}_{2r\tau+k}^i - \eta \left( \mathbf{g}_{2r\tau+k}^i - \hat{\mathbf{g}}_{(r\tau)}^i + \hat{\mathbf{g}}_{(r\tau)} \right)$ 
12:    end for
13:  end for
14:  Compute:  $\bar{\mathbf{x}}_{2(r+1)\tau} = \bar{\mathbf{x}}_{2r\tau} - \frac{\eta}{n} \sum_{j=1}^n \sum_{l=\tau}^{2\tau-1} \mathbf{g}_{2r\tau+l}^j$ 
15:  Broadcast:  $\mathbf{x}_{2(r+1)\tau}^i = \bar{\mathbf{x}}_{2(r+1)\tau}$ , for each  $i \in [n]$ 
16: end for
```

Assumptions (gradient similarity)

Assumption 1 (Standard gradient similarity – SGS)

For some $\zeta \geq 0$, we have

$$\sup_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|_2^2 \leq \zeta^2.$$

Assumption 1+ (Uniform gradient similarity – UGS)

For some $\bar{\zeta} \geq 0$, we have

$$\sup_{\mathbf{x} \in \mathbb{R}^d} \sup_{i \in [n]} \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|_2^2 \leq \bar{\zeta}^2$$

Assumptions (Hessian similarity)

Assumption 2 (Standard Hessian similarity – SHS)

For some $\delta \in [0, 2L]$, we have

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x}) - \nabla f_i(\mathbf{y}) + \nabla f(\mathbf{y})\|_2^2 \leq \delta^2 \|\mathbf{x} - \mathbf{y}\|_2^2,$$

for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.

Assumption 2+ (Uniform Hessian similarity – UHS)

For some $\bar{\delta} \in [0, 2L]$, we have

$$\|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x}) - \nabla f_i(\mathbf{y}) + \nabla f(\mathbf{y})\|_2 \leq \bar{\delta} \|\mathbf{x} - \mathbf{y}\|_2,$$

for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and for all $i \in [n]$.

Assumptions (weak convexity, Lipschitz continuous Hessian)

Assumption 3 (Weak convexity – WC)

For some $\rho \in [0, L]$, we have

$$f_i(\mathbf{x}) + \frac{\rho}{2} \mathbf{x}^\top \mathbf{x} \text{ is convex,}$$

for all $i \in [n]$.

Assumption 4 (Lipschitz continuous Hessian – LCH)

For some $\mathcal{M} \geq 0$, there exists (at least) one function \hat{f} such that:
 $\hat{f} \in \mathbf{conv}\{f_1, \dots, f_n\}$, and

$$\left\| \nabla^2 \hat{f}(\mathbf{x}) - \nabla^2 \hat{f}(\mathbf{y}) \right\|_2 \leq \mathcal{M} \|\mathbf{x} - \mathbf{y}\|_2,$$

for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.

Lemma 1

There exists $\eta > 0$ such that MbSGD ensures the following upper bound on $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|_2^2$:

$$\mathcal{O} \left(\frac{L\Delta}{R} + \sqrt{\frac{L\Delta\sigma^2}{n\tau R}} \right).$$

LocalSGD: non-convex speedup from WC

Lemma 2 ([Kol+20])

Under Assumption 1, there exists $\eta > 0$ such that LocalSGD ensures the following upper bound on $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|_2^2$:

$$\mathcal{O} \left(\frac{L\Delta}{R} + \sqrt{\frac{L\Delta\sigma^2}{n\tau R}} + \left(\frac{L\Delta\zeta}{R} \right)^{\frac{2}{3}} + \frac{(L\Delta\sigma)^{\frac{2}{3}}}{\tau^{\frac{1}{3}} R^{\frac{2}{3}}} \right).$$

Theorem 1 (Ours)

Under Assumptions 1 and 3, there exists $\eta > 0$ such that LocalSGD ensures the following upper bound on $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|_2^2$:

$$\mathcal{O} \left(\left(\frac{L}{\tau} + \rho \right) \frac{\Delta}{R} + \sqrt{\frac{L\Delta\sigma^2}{n\tau R}} + \left(\frac{L\Delta\zeta}{R} \right)^{\frac{2}{3}} + \frac{(L\Delta\sigma)^{\frac{2}{3}}}{\tau^{\frac{1}{3}} R^{\frac{2}{3}}} \right).$$

LocalSGD: convex speedup without UGS

Lemma 3 ([WPS20])

Under Assumption 1+, if all f_i are convex, $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$, and $\|\bar{\mathbf{x}}_0 - \mathbf{x}^*\|_2 \leq D$, then there exists $\eta > 0$ such that LocalSGD ensures the following upper bound on $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[f(\bar{\mathbf{x}}_t)] - f^*$:

$$\mathcal{O} \left(\frac{LD^2}{\tau R} + \frac{\sigma D}{\sqrt{n\tau R}} + \left(\frac{L\bar{\zeta}^2 D^4}{R^2} \right)^{\frac{1}{3}} + \left(\frac{L\sigma^2 D^4}{\tau R^2} \right)^{\frac{1}{3}} \right).$$

Theorem 2 (Ours)

Under Assumption 1, ...

$$\mathcal{O} \left(\frac{LD^2}{\tau R} + \frac{\sigma D}{\sqrt{n\tau R}} + \left(\frac{L\zeta^2 D^4}{R^2} \right)^{\frac{1}{3}} + \left(\frac{L\sigma^2 D^4}{\tau R^2} \right)^{\frac{1}{3}} \right).$$

Theorem 3 (Ours)

Under Assumptions 1, 2+ and 4, there exists $\eta > 0$ such that LocalSGD ensures the following upper bound on $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|_2^2$:

$$\mathcal{O}\left(\frac{L\Delta}{R} + \sqrt{\frac{L\Delta\sigma^2}{n\tau R}} + \left(\frac{\bar{\delta}\Delta\zeta}{R}\right)^{\frac{2}{3}} + \frac{(L\Delta\sigma)^{\frac{2}{3}}}{\tau^{\frac{1}{3}}R^{\frac{2}{3}}} + \left(\frac{\mathcal{M}^2\Delta^4\zeta^4}{R^4}\right)^{\frac{1}{5}}\right).$$

SCAFFOLD: existing analyses

Lemma 4 ([Kar+20])

Suppose in Line 14 of Algorithm 1, a different global stepsize η_g can be used when aggregating the updates. There exists $\eta_g \geq \eta > 0$ such that SCAFFOLD ensures the following upper bound on $\frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_{2r\tau})\|_2^2$:

$$\mathcal{O} \left(\frac{L\Delta}{R} + \sqrt{\frac{L\Delta\sigma^2}{n\tau R}} \right).$$

Lemma 5 ([Kar+20])

Suppose $\hat{\mathbf{g}}_{(r\tau)}^i = \nabla f_i(\bar{\mathbf{x}}_{2r\tau})$ in Line 6 of Algorithm 1. Under Assumptions 2+ and 3, if all f_i are quadratic, then there exists $\eta > 0$ such that SCAFFOLD ensures the following upper bound on $\frac{2}{T} \sum_{r=0}^{R-1} \sum_{k=0}^{\tau-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_{2r\tau+\tau+k})\|_2^2$:

$$\mathcal{O} \left(\left(\frac{L}{\tau} + \bar{\delta} + \rho \right) \frac{\Delta}{R} + \sqrt{\frac{L\Delta\sigma^2}{n\tau R}} \right).$$

SCAFFOLD: speedup without quadratic/UHS

Theorem 4 (Ours)

Under Assumptions 2 and 3, there exists $\eta > 0$ such that SCAFFOLD ensures the following upper bound on $\frac{2}{T} \sum_{r=0}^{R-1} \sum_{k=0}^{\tau-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_{2r\tau+\tau+k})\|_2^2$:

$$\mathcal{O} \left(\left(\frac{L}{\tau} + \sqrt{L\delta} + \rho \right) \frac{\Delta}{R} + \sqrt{\frac{L\Delta\sigma^2}{n\tau R}} + \frac{(L\Delta\sigma)^{\frac{2}{3}}}{\tau^{\frac{1}{3}} R^{\frac{2}{3}}} \right).$$

Theorem 5 (Ours)

Under Assumptions 2 to 4 with $\mathcal{M} = 0$, there exists $\eta > 0$ s.t. SCAFFOLD ensures the following upper bound on $\frac{2}{T} \sum_{r=0}^{R-1} \sum_{k=0}^{\tau-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_{2r\tau+\tau+k})\|_2^2$:

$$\mathcal{O} \left(\left(\frac{L}{\tau} + \sqrt{\bar{\delta}\delta} + \rho \right) \frac{\Delta}{R} + \sqrt{\frac{L\Delta\sigma^2}{n\tau R}} + \frac{(\bar{\delta}\Delta\sigma)^{\frac{2}{3}}}{\tau^{\frac{1}{3}} R^{\frac{2}{3}}} \right).$$

Proof for SCAFFOLD: descent lemma

Notations. For $t = r\tau + k$, $r \in [0, R - 1]$, $k \in [0, \tau - 1]$, we denote

$$r(t) = r\tau, \quad \mathbf{x}_{(t)}^i = \mathbf{x}_{2r\tau + \tau + k}^i, \quad \mathbf{g}_{(t)}^i = \mathbf{g}_{2r\tau + \tau + k}^i,$$

and

$$\bar{\mathbf{x}}_{(t)}^i = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{(t)}^i, \quad \Xi_{(t)} = \frac{1}{n} \sum_{i=1}^n \left\| \mathbf{x}_{(t)}^i - \bar{\mathbf{x}}_{(t)} \right\|_2^2.$$

Lemma 6

For $\eta \leq \frac{1}{2L}$, SCAFFOLD ensures

$$\begin{aligned} & \frac{2}{T} \sum_{t=0}^{R\tau-1} \left[\mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}_{(t)}) \right\|_2^2 + \frac{1}{2} \mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_{(t)}^i) \right\|_2^2 \right] \\ & \leq \frac{4\Delta}{\eta T} + 2\eta \frac{L\sigma^2}{n} + \frac{2}{T} \sum_{r=0}^{R-1} \sum_{k=0}^{\tau-1} L^2 \mathbb{E} [\Xi_{(r\tau+k)}]. \end{aligned}$$

Proof for SCAFFOLD: distance lemma

Lemma 7

Under Assumptions 2 and 3, for $\gamma = \frac{1}{3(\tau-1)}$ and $\eta \leq \frac{(1-\rho/L)}{2\rho}\gamma$, we have

$$\begin{aligned} \mathbb{E} [\Xi(t)] &\leq 3\gamma^{-2}(1 + \gamma^{-1})\eta^4\delta^2 \sum_{l=0}^{k-1} \mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_{(r(t)+l)}^i) \right\|_2^2 \\ &\quad + 3 \left(\frac{1 + \gamma^{-1}}{\tau} + 1 \right) \eta^2 k \sigma^2 + 3\gamma^{-2}\eta^4 k \delta^2 \frac{\sigma^2}{n}, \end{aligned} \quad (1)$$

where $k = t - r(t)$.

References

- [Kar+20] Sai Praneeth Karimireddy et al. “Scaffold: Stochastic controlled averaging for federated learning”. In: *International Conference on Machine Learning*. PMLR. 2020, pp. 5132–5143.
- [Kol+20] Anastasia Koloskova et al. “A unified theory of decentralized sgd with changing topology and local updates”. In: *International Conference on Machine Learning*. PMLR. 2020, pp. 5381–5393.
- [WPS20] Blake E Woodworth, Kumar Kshitij Patel, and Nati Srebro. “Minibatch vs local sgd for heterogeneous distributed learning”. In: *Advances in Neural Information Processing Systems 33* (2020), pp. 6281–6292.