

# CS331-HW6-Lukang-Sun

October 16, 2021

---

**Algorithm 1** DCGD-SHIFT

---

**Parameters:** shift  $h_i$ , learning rate  $\gamma > 0$ , starting point  $x^0 \in \mathbb{R}^d$ , compression operators  $\mathcal{C}_1 \in \mathbb{B}^d(\omega_1), \dots, \mathcal{C}_n \in \mathbb{B}^d(\omega_n), \mathcal{C} = Id$

- 1: **for**  $k = 0, 1, 2, \dots$  **do**
- 2:   **for** all workers  $i \in \{1, 2, \dots, n\}$  in parallel **do**
- 3:     Compute local gradient  $\nabla f_i(x^k)$
- 4:     Compress local gradient  $g_i^k = \mathcal{C}_i^k(\nabla f_i(x^k) - h_i)$
- 5:     Send  $g_i^k$  to master
- 6:   **end for**
- 7:   Master computes the aggregate  $\hat{g}^k = \frac{1}{n} \sum_{i=1}^n g_i^k + \frac{1}{n} \sum_{i=1}^n h_i$
- 8:   Master broadcasts the compressed aggregate  $g^k = \mathcal{C}(\hat{g}^k)$  to all workers
- 9:   **for** all workers  $i \in \{1, 2, \dots, n\}$  in parallel **do**
- 10:     Compute the next iterate  $x^{k+1} = \text{prox}_{\gamma R}(x^k - \gamma g^k)$
- 11:   **end for**
- 12: **end for**

---

**p1.**

**Theorem.** Assume  $f_i$  is convex and  $L_i$ -smooth for all  $i$ , and let  $f$  be  $L$ -smooth. Let the gradient estimator  $g$  be defined as in Algorithm 1, where

$$\mathcal{C}_1 \in \mathbb{B}^d(\omega_1), \quad \mathcal{C}_2 \in \mathbb{B}^d(\omega_2), \quad \dots, \quad \mathcal{C}_n \in \mathbb{B}^d(\omega_n), \quad \mathcal{C} = Id$$

are independent compression operators,  $h_i = \nabla f_i(y)$ . Then

$$G(x, y) \leq 2AD_f(x, y)$$

where

$$A = L + \max_i \left( L_i \frac{\omega_i}{n} \right).$$

Let the step size  $0 \leq \gamma \leq \frac{1}{A}$ , then

$$\mathbb{E} \left[ \|x^k - x^*\|^2 \right] \leq (1 - \gamma\mu)^k \|x^0 - x^*\|^2. \quad (1)$$

*Proof.*  $\mathcal{C}_i \in \mathbb{B}^d(\omega_i)$  are independent and  $\mathcal{C}(x) = x$ . We get the estimate

$$\begin{aligned} \mathbb{E} [\|g(x) - \nabla f(y)\|^2] &= \mathbb{E} [\|\hat{g}(x) - \nabla f(x)\|^2] + \mathbb{E} [\|\mathbb{E}[\hat{g}(x)] - \nabla f(y)\|^2] \\ &\leq \mathbb{E} [\|\hat{g}(x) - \nabla f(x)\|^2] + 2LD_f(x, y) \end{aligned}$$

We now bound each of the above terms individually. Let  $h_i = \nabla f_i(y)$ ,  $a_i \stackrel{\text{def}}{=} g_i(x) + h_i - \nabla f_i(x)$  and note that  $\mathbb{E}[a_i] = 0$ . The second term can be bounded as

$$\begin{aligned} \mathbb{E} [\|\hat{g}(x) - \nabla f(x)\|^2] &= \mathbb{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n (g_i(x) + h_i - \nabla f_i(x)) \right\|^2 \right] \\ &= \mathbb{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n a_i \right\|^2 \right] \\ &= \frac{1}{n^2} \mathbb{E} \left[ \sum_{i=1}^n \|a_i\|^2 + \sum_{i \neq j} \langle a_i, a_j \rangle \right] \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [\|a_i\|^2] + \sum_{i \neq j} \mathbb{E} [\langle a_i, a_j \rangle] \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [\|a_i\|^2] + \sum_{i \neq j} \underbrace{\mathbb{E}[a_i]}_{=0} \underbrace{\mathbb{E}[a_j]}_0 \\ &\leq \frac{1}{n^2} \sum_{i=1}^n \omega_i \|\nabla f_i(x) - \nabla f_i(y)\|^2 \\ &\leq \frac{1}{n^2} \sum_{i=1}^n 2\omega_i L_i D_{f_i}(x, y) \\ &\leq \frac{2 \max_i \omega_i L_i}{n} \frac{1}{n} \sum_{i=1}^n D_{f_i}(x, y) = \frac{2 \max_i \omega_i L_i}{n} D_f(x, y), \end{aligned}$$

so finally

$$\mathbb{E} [\|g(x) - \nabla f(y)\|^2] \leq 2(L + \frac{\max_i \omega_i L_i}{n}) D_f(x, y)$$

Inequality (1) is the corollary under AC-condition,  $A = L + \max_i (L_i \frac{\omega_i}{n})$ ,  $C = 0$ .  $\square$

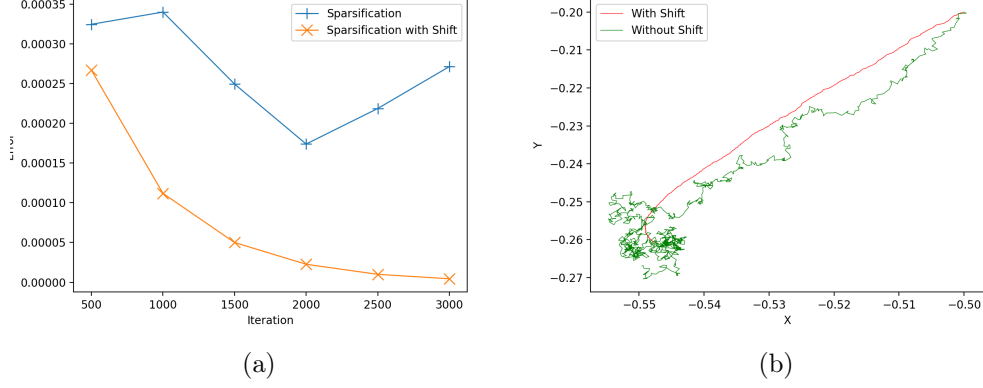


Figure 1: (a) shows the error in terms of iteration, (b) shows the trajectories of SGD-IND with shift and SGD-IND without shift

**p2.** In my experiments (see Figure 1.), I set  $d = 2, n = 10, f(x) = \frac{1}{10} \sum_{i=1}^{10} f_i(x), f_i(x) = \frac{1}{2} \|a_i x - b_i\|_2^2, a = [\text{matrix}([0.94884523, 0.31257516]), \text{matrix}([0.64695759, 0.79089169]), \text{matrix}([0.70218109, 0.91473775]), \text{matrix}([0.03035042, 0.21034799]), \text{matrix}([0.99278455, 0.2554682]), \text{matrix}([0.16064759, 0.09062056]), \text{matrix}([0.41438167, 0.77718962]), \text{matrix}([0.4953842, 0.93027311]), \text{matrix}([0.7692516, 0.19772597]), \text{matrix}([0.12430258, 0.03779965])], b = [\text{matrix}([0.4231786]), \text{matrix}([0.524466]), \text{matrix}([0.17981303]), \text{matrix}([0.50033184]), \text{matrix}([0.71116473]), \text{matrix}([0.0604399]), \text{matrix}([0.37100656]), \text{matrix}([0.91263449]), \text{matrix}([0.66417151]), \text{matrix}([0.65995097])]$ , based on these information, we can get that  $x_\star = \text{matrix}([-0.54051746] [-0.26890662])$ , I sample 10 points in the last step to estimate  $E[\|x_{\text{last step}} - x_\star\|^2]$ . I use SGD-IND ( $p = 0.5$ ) with and without shift in my experiments, the experiments show that SGD-IND without shift will converge to a neighborhood of the optimal point, but SGD-IND with shift will converge to the optimal point, these results quite match the theory.

**p3.** (1.) Since in this case

$$g^k = g(x^k) - g(y^k) + \nabla f(y^k) = \nabla f(x^k),$$

which is exactly the descent direction for GD. (2.) I don't find Corollary 95. Do you mean Corollary 97?

**Theorem.** Assume  $f$  is  $\mu$ -convex and  $L$ -smooth. The gradient estimator  $g = \nabla f + \xi$  is unbiased and satisfies the expected smoothness bound

$$\sigma(x) \stackrel{\text{def}}{=} E[\|g(x) - g(x^\star)\|^2] \leq 2LD_f(x, x^\star)$$

Then  $L$ -SVRG with stepsize  $\gamma = \frac{1}{6L}$  satisfies

$$\mathbb{E} [V^k] \leq \left(1 - \min \left\{ \frac{\mu}{6L}, \frac{p}{2} \right\}\right)^k V^0$$

where

$$V^k \stackrel{(183)}{=} \|x^k - x^\star\|^2 + M\gamma^2\sigma^k, \quad M = \frac{4}{p}, \quad \sigma^k = \sigma(y^k) \stackrel{def}{=} \mathbb{E} \left[ \|g(y^k) - g(x^\star)\|^2 \right]$$

So,

$$k \geq \max \left\{ \frac{6L}{\mu}, \frac{2}{p} \right\} \log \frac{1}{\varepsilon} \Rightarrow \mathbb{E} [V^k] \leq \varepsilon V^0$$

since  $\max \left\{ \frac{6L}{\mu}, \frac{2}{p} \right\} > \frac{L}{\mu}$ , so this kind of analysis is worse than GD. When  $p$  is not small, the order of convergence rate is the same between GD and  $L$ -SVRG( $g = \nabla f + \xi$ )(ignore the constant factor), when  $p$  is small (like  $p = \epsilon$ ), then the method in this problem is much worse than GD.

**p4.** The expected smoothness constant is

$$A'' = \frac{n - \tau}{\tau(n - 1)} \max_i L_i + \frac{n(\tau - 1)}{\tau(n - 1)} L$$

In each iteration, one has to compute  $2|S| = 2\tau$  1-gradient and  $p$   $n$ -gradient, so

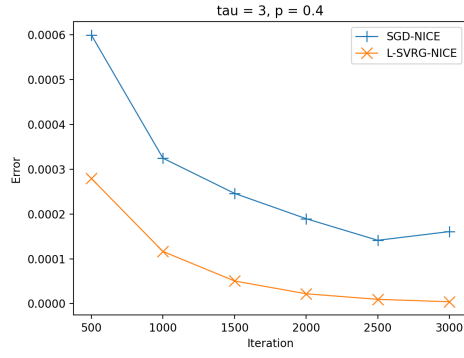
$$\mathbf{COST} = 2\tau + pn,$$

so the total complexity is

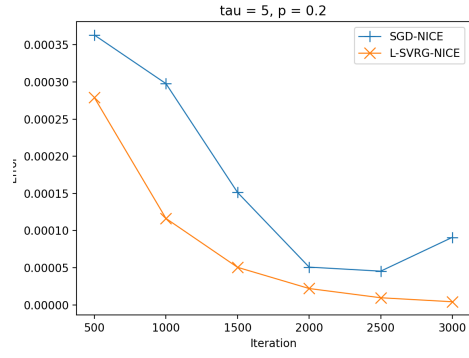
$$\max \left\{ \frac{6\left(\frac{n-\tau}{\tau(n-1)} \max_i L_i + \frac{n(\tau-1)}{\tau(n-1)} L\right)(2\tau + pn)}{\mu}, \frac{2(2\tau + pn)}{p} \right\} \log\left(\frac{1}{\epsilon}\right),$$

Computing the minimum point of  $\left(\frac{n-\tau}{\tau(n-1)} \max_i L_i + \frac{n(\tau-1)}{\tau(n-1)} L\right)(2\tau + pn)$ , it is  $\tau_\star = \sqrt{\frac{(\max L_i - L)n^2 p}{2(nL - \max L_i)}}$ , since optimal  $\tau$  should be an integer between 1 and  $n-1$ , we find optimal  $\tau^\star$  between  $\lfloor \tau_\star \rfloor$  and  $\lceil \tau_\star \rceil$  (only consider  $\lfloor \tau_\star \rfloor, \lceil \tau_\star \rceil \in [1, n-1]$ , if one of them is not in  $[1, n-1]$ , ignore it), consider the  $\hat{\tau}$ , such that  $\frac{6A''}{\mu} = \frac{1}{p}$ , if  $\hat{\tau} > \tau^\star$ , choose  $\tau_{min} = \tau^\star$ , else  $\tau_{min} = \text{optimal}\{\lfloor \hat{\tau} \rfloor, \lceil \hat{\tau} \rceil \in [1, n-1]\}$ .  $\tau_{min}$  depends on  $p$ , if  $p$  is small, then  $\tau_{min}$  is also small, if  $p$  is big, then  $\tau_{min}$  will also get bigger.

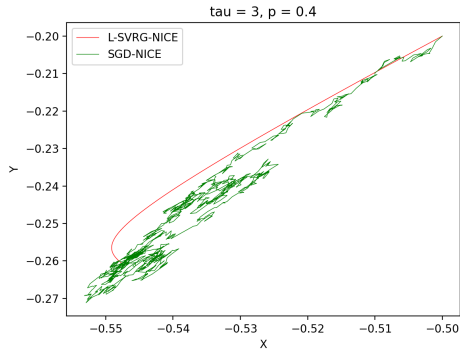
**p5.** In my experiments(see Figure2.), I set  $d = 2, n = 10, f(x) = \frac{1}{10} \sum_{i=1}^{10} f_i(x), f_i(x) = \frac{1}{2} \|a_i x - b_i\|_2^2, a = [\text{matrix}([0.94884523, 0.31257516]), \text{matrix}([0.64695759, 0.79089169]), \text{matrix}([0.70218109, 0.91473775]), \text{matrix}([0.03035042, 0.21034799]), \text{matrix}([0.99278455, 0.2554682]), \text{matrix}([0.16064759, 0.09062056]), \text{matrix}([0.41438167, 0.77718962]), \text{matrix}([0.4953842, 0.93027311]), \text{matrix}([0.7692516, 0.19772597]), \text{matrix}([0.12430258, 0.03779965])], b = [\text{matrix}([0.4231786]), \text{matrix}([0.524466]), \text{matrix}([0.17981303]), \text{matrix}([0.50033184]), \text{matrix}([0.71116473]), \text{matrix}([0.0604399]), \text{matrix}([0.37100656]), \text{matrix}([0.91263449]), \text{matrix}([0.66417151]), \text{matrix}([0.65995097])]$ , based on these information, we can get that  $x_\star = \text{matrix}([-0.54051746] [-0.26890662])$ , I sample 10 points in the last step to estimate  $E[|x_{\text{last step}} - x_\star|^2]$ . The theory states that L-SVRG-NICE will converge to the optimal point while SGD-NICE will only converge to the neighborhood of the optimal point, this quite match the experiments results. You can see that the error of L-SVRG-NICE keeps dropping while SGD-NICE not and the trajectory of L-SVRG-NICE converge to the optimal point smoothly but the trajectory of SGD-NICE converges to the neighborhood of the optimal point with fluctuation.



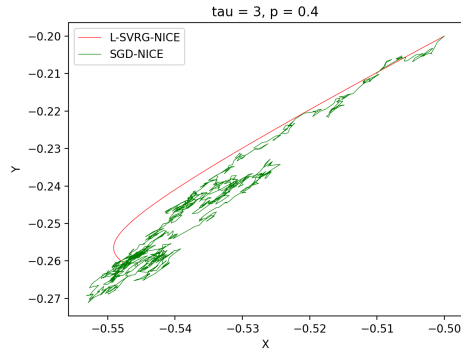
(a)



(b)



(c)



(d)

Figure 2: (a) shows the error in terms of iteration with  $\tau = 3, p = 0.4$  (b) shows the error in terms of iteration with  $\tau = 5, p = 0.2$ , (c) shows the trajectories of SGD-NICE and L-SVRG-NICE with  $\tau = 3, p = 0.4$ , (d) shows the trajectories of SGD-NICE and L-SVRG-NICE with  $\tau = 5, p = 0.2$